

## PATENT ABSTRACTS OF JAPAN

(11)Publication number : 10-272000  
 (43)Date of publication of application : 13.10.1998

(51)Int. Cl. C12Q 1/68  
 C12N 15/09  
 G06F 17/30  
 // G01N 33/50

(21)Application number : 09-369833 (71)Applicant : AFFYMETRIX INC  
 (22)Date of filing : 11.12.1997 (72)Inventor : WEBSTER TERESA A  
 MORRIS MACDONALD S  
 MITTMANN MICHAEL P  
 LOCKHART DAVID J  
 HO MING-HSIU  
 BERNHART DEREK  
 JEVONS LUIS C

## (30)Priority

Priority number : 96 33053 Priority date : 12.12.1996 Priority country : US  
 97 828952 28.03.1997

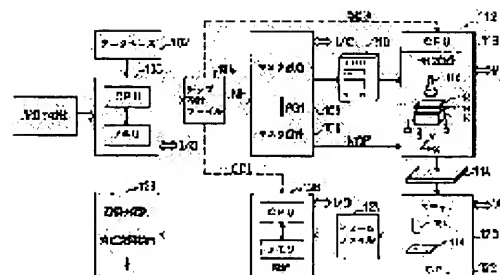
US

## (54) ANALYTIC TECHNIQUE OF BIOLOGICAL SEQUENCE UTILIZING COMPUTER

## (57)Abstract:

PROBLEM TO BE SOLVED: To analyze a biological sequence of nucleic acid, etc., with computer system by inputting and analyzing plural base calls to each base position along a part of nucleic acid sequence of a specimen and producing and displaying a single base call.

SOLUTION: In a method for analyzing a nucleic acid sequence of a specimen in a computer system, plural base calls to each base position along at least a part are inputted to a computer 100 and a base call having highest occurrence frequency among these plural base calls is obtained by database 102 and a single base call is produced as a base call having highest frequency of occurrence and inputted to a chip design file 104 and a pattern is arranged on a mask 110 and the chip design file 104 and an image file 124 are analyzed and single base call each led from plural base calls to a specific base position in the nucleic acid sequence of the specimen is displayed as an output 128.



## LEGAL STATUS

[Date of request for examination] 22.05.2002

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application]

converted registration]

[Date of final disposal for application]

[Patent number]

[Date of registration]

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's decision of rejection]

[Date of extinction of right]

Copyright (C); 1998, 2003 Japan Patent Office

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平10-272000

(43) 公開日 平成10年(1998)10月13日

(51) Int. Cl. <sup>6</sup>	識別記号	P I
C 1 2 Q 1/68		C 1 2 Q 1/68 A
C 1 2 N 15/00		G 0 1 N 33/50 P
G 0 6 F 17/30		C 1 2 N 15/00 A
// G 0 1 N 33/50		G 0 6 F 15/40 3 7 0 F

審査請求 未請求 請求項の数35 F D 外国語出願 (全104頁)

(21) 出願番号 特願平9-369833

(22) 出願日 平成9年(1997)12月11日

(31) 優先権主張番号 60/033, 053

(32) 優先日 1996年12月12日

(33) 優先権主張国 米国 (US)

(31) 優先権主張番号 08/828, 952

(32) 優先日 1997年3月28日

(33) 優先権主張国 米国 (US)

(71) 出願人 598012740

アフィメトリックス・インコーポレーテッド

APFYMETRIX INCORPORATED

アメリカ合衆国 カリフォルニア州95051  
サンタ・クララ、セントラル・エクスプレ  
スウェイ, 3380

(72) 発明者 テレサ・エー・ウェブスター

アメリカ合衆国 カリフォルニア州94021  
ローマ・マー, ベスカドーレ・ロード,  
10002

(74) 代理人 弁理士 五十嵐 孝雄 (外2名)

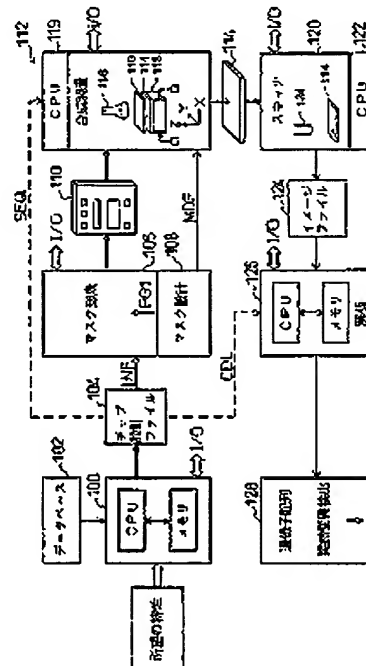
最終頁に続く

(54) 【発明の名称】 コンピュータを利用した生物学的配列の解析技術

(57) 【要約】

【課題】 核酸のような生物学的配列を解析するためのコンピュータ支援技術を提供する。

【解決手段】 コンピュータシステムは、核酸プローブと検体核酸配列との間のハイブリッド形成親和性を示すハイブリッド形成強度を解析して、検体配列における塩基を呼び出す（塩基コールを行う）。複数の塩基コールを組み合わせることで、単一の塩基コールを形成する。また、コンピュータシステムは、ハイブリッド形成強度を解析することによって、遺伝子発現を監視したり、あるいは、基盤からの遺伝子発現の変化を監視したりする。



## 【特許請求の範囲】

【請求項1】 コンピュータシステムにおいて、検体の核酸配列を解析するための方法であって、前記検体の核酸配列の少なくとも一部分に沿った各塩基位置に対して、複数の塩基コールを入力する工程と、各塩基位置に対して、前記複数の塩基コールを解析して単一塩基コールを生成する工程と、前記検体の核酸配列の少なくとも一部分に沿った複数の塩基位置に対する複数の単一塩基コールであって、前記検体の核酸配列における特定の塩基位置に対する前記複数の塩基コールからそれぞれ導かれた単一塩基コールを表示する工程と、を備える方法。

【請求項2】 請求項1記載の方法であって、前記解析工程は、各塩基位置に対して、前記複数の塩基コールの中から最も発生頻度の高い塩基コールを求める工程と、前記単一塩基コールを、前記塩基位置において最も発生頻度の高い塩基コールとして生成する工程と、を備える方法。

【請求項3】 請求項1記載の方法であって、更に、ユーザーによる起動に応じて各塩基位置に前記複数の塩基コールを表示させるためのスクリーンアイコンを表示する工程、を備える方法。

【請求項4】 請求項1記載の方法であって、更に、ユーザーによる起動に応じて各塩基位置に前記複数の塩基コールを表示させないためのスクリーンアイコンを表示する工程、を備える方法。

【請求項5】 請求項1記載の方法であって、更に、塩基位置に従って、前記単一塩基コールと並べて各塩基における前記複数の塩基コールを表示する工程、を備える方法。

【請求項6】 請求項5記載の方法であって、更に、前記複数の塩基コールの各々と共に、プローブと前記検体の核酸配列とのハイブリッド形成親和性を示すハイブリッド形成強度を表示する工程を備え、各塩基コールは前記ハイブリッド形成強度の解析により求められる、方法。

【請求項7】 コンピュータシステムにおいて、検体の核酸配列における未知の塩基を呼び出すための方法であって、複数の核酸プローブに関するハイブリッド形成強度であって、それぞれ核酸プローブと前記検体の核酸配列との間のハイブリッド形成親和性を示すハイブリッド形成強度を入力する工程と、各組のプローブに関して前記未知の塩基に対する塩基コールを演算する工程と、前記複数の組のプローブに関して最も発生頻度の高い前記未知の塩基に対する塩基コールに従って、前記複数の組のプローブに対する単一塩基コールを演算する工程と、を備える方法。

【請求項8】 請求項7記載の方法であって、各組のプローブが同一の参照配列に従って生成されたものである、方法。

【請求項9】 請求項7記載の方法であって、更に、所定の条件下で、前記複数の核酸プローブに関する前記単一塩基コールを特定する例外規則をチェックする工程、を備える方法。

【請求項10】 コンピュータシステムにおいて、コンピュータによる塩基呼び出し処理に関するパラメータを動的に変更する方法であって、ユーザーによって変更可能なパラメータを含む前記塩基呼び出し処理を利用して、検体の核酸配列の少なくとも一部に関する塩基コールを生成する工程と、前記検体の核酸配列の少なくとも一部に関する複数の塩基コールを表示する工程と、前記塩基呼び出し処理のパラメータを表示する工程と、前記塩基呼び出し処理のパラメータに対する新しい値を規定する入力をユーザーから受け取る工程と、前記塩基呼び出し処理と前記パラメータに関する新しい値とを用いて、前記検体の核酸配列の少なくとも一部に関して更新された塩基コールを生成する工程と、前記検体の核酸配列の少なくとも一部に関して前記更新された塩基コールを表示する工程と、を備える方法。

【請求項11】 請求項10記載の方法であって、更に、前記塩基呼び出し処理のための、複数のユーザー変更可能なパラメータを表示する工程、を備える方法。

【請求項12】 請求項10記載の方法であって、前記パラメータが、定数、閾値、および、範囲から成るグループから選択される、方法。

【請求項13】 コンピュータシステムにおいて、検体の核酸配列における遺伝子の発現を監視する方法であって、前記遺伝子に対して完全に相補的である完全対合プローブと前記遺伝子に対して不対合塩基を少なくとも一つ備える不対合プローブとの複数のプローブ対に関する複数のハイブリッド形成強度であって、前記完全対合および不対合プローブと前記検体の核酸配列との間のハイブリッド形成親和性を示すハイブリッド形成強度を入力する工程と、前記検体の核酸配列の遺伝子発現コールを生成するために、各対の完全対合プローブと不対合プローブのハイブリッド形成強度を比較する工程と、前記遺伝子発現コールを表示する工程と、を備える方法。

【請求項14】 請求項13記載の方法であって、更に、或る塩基位置における完全対合プローブと不対合プローブのハイブリッド形成強度の差分を、差分の閾値と比較する工程、を備える方法。



【請求項15】 請求項13記載の方法であって、更に、  
或る塩基位置における完全対合ブローブと不対合ブローブのハイブリッド形成強度の比を、比の閾値と比較する工程、を備える方法。

【請求項16】 請求項13記載の方法であって、更に、  
決定行列を用いて前記遺伝子発現コールを決定する工程、を備える方法。

【請求項17】 請求項13記載の方法であって、  
前記遺伝子発現コールが、発現している、どちらともいえない、発現していない、から成るグループから選択される、方法。

【請求項18】 コンピュータシステムにおいて、検体の核酸配列における遺伝子の発現を監視する方法であって、

前記遺伝子に対して完全に相補的である完全対合ブローブと前記遺伝子に対して不対合な塩基を少なくとも一つ備える不対合ブローブとの複数のブローブ対に関する複数のハイブリッド形成強度であって、前記完全対合および不対合ブローブと前記検体の核酸配列との間のハイブリッド形成親和性を示すハイブリッド形成強度を入力する工程と、  
各対の完全対合ブローブと不対合ブローブのハイブリッド形成強度を比較する工程と、  
前記検体の核酸配列の遺伝子発現コールを生成する工程と、を備える方法。

【請求項19】 請求項18記載の方法であって、更に、  
或る塩基位置における完全対合ブローブと不対合ブローブのハイブリッド形成強度の差分を、差分の閾値と比較する工程、を備える方法。

【請求項20】 請求項18記載の方法であって、更に、  
或る塩基位置における完全対合ブローブと不対合ブローブのハイブリッド形成強度の比を、比の閾値と比較する工程、を備える方法。

【請求項21】 請求項18記載の方法であって、更に、  
決定行列を用いて前記遺伝子発現コールを求める工程、を備える方法。

【請求項22】 請求項18記載の方法であって、  
前記遺伝子発現コールが、発現している、どちらともいえない、発現していない、から成るグループから選択される、方法。

【請求項23】 コンピュータシステムにおいて、検体の核酸配列における遺伝子の発現変化を監視する方法であって、  
前記遺伝子に対して完全に相補的である完全対合ブローブと前記遺伝子に対して不対合な塩基を少なくとも一つ

備える不対合ブローブとの複数のブローブ対に関する複数のハイブリッド形成強度であって、前記完全対合および不対合ブローブと前記検体の核酸配列との間のハイブリッド形成親和性を示すハイブリッド形成強度を入力する工程と、

前記検体の核酸配列の遺伝子発現レベルを生成するために、各対の完全対合ブローブと不対合ブローブのハイブリッド形成強度を比較する工程と、

前記遺伝子発現レベルを、基準の遺伝子発現レベルと比較することにより発現変化を求める工程と、

前記検体の核酸配列における前記遺伝子の発現変化を表示する工程と、を備える方法。

【請求項24】 請求項23記載の方法であって、  
前記発現変化がグラフで表示される、方法。

【請求項25】 請求項23記載の方法であって、更に、  
請求項23の前記入力工程と前記比較工程に従って、前記基準の遺伝子発現レベルを生成する工程、を備える方法。

【請求項26】 請求項23記載の方法であって、更に、  
前記検体の核酸配列とハイブリッドを形成する完全対合および不対合ブローブのハイブリッド形成強度と、基準配列とハイブリッドを形成する完全対合および不対合ブローブのハイブリッド形成強度とを、差分の閾値に対して比較する工程、を備える方法。

【請求項27】 請求項23記載の方法であって、更に、  
前記検体の核酸配列とハイブリッドを形成する完全対合および不対合ブローブのハイブリッド形成強度と、基準配列とハイブリッドを形成する完全対合および不対合ブローブのハイブリッド形成強度とを、比の閾値に対して比較する工程、を備える方法。

【請求項28】 請求項23記載の方法であって、更に、  
決定行列を用いて前記検体の核酸配列における前記遺伝子の発現変化を決定する工程、を備える方法。

【請求項29】 請求項23記載の方法であって、  
前記検体の核酸配列における前記遺伝子の発現変化が、増加、増加傾向、減少、減少傾向、変化なし、から成るグループから選択される、方法。

【請求項30】 コンピュータシステムにおいて、検体の核酸配列における遺伝子の発現変化を監視する方法であって、

前記遺伝子に対して完全に相補的である完全対合ブローブと前記遺伝子に対して不対合な塩基を少なくとも一つ備える不対合ブローブとの複数のブローブ対に関する複数のハイブリッド形成強度であって、前記完全対合および不対合ブローブと前記検体の核酸配列との間のハイブリッド形成親和性を示すハイブリッド形成強度を入力す

10

20

30

40

50

る工程と、  
前記検体の核酸配列の遺伝子発現レベルを生成するために、各対の完全対合ブローブと不対合ブローブのハイブリッド形成強度を比較する工程と、  
前記遺伝子発現レベルを、基準の遺伝子発現レベルと比較することにより発現変化を求める工程と、を備える方法。

【請求項31】 請求項30記載の方法であって、更に、  
請求項30の前記入力工程と前記比較工程に従って、前記基準の遺伝子発現レベルを生成する工程、を備える方法。

【請求項32】 請求項30記載の方法であって、更に、  
前記検体の核酸配列とハイブリッドを形成する完全対合および不対合ブローブのハイブリッド形成強度と、基準配列とハイブリッドを形成する完全対合および不対合ブローブのハイブリッド形成強度とを、差分の閾値に対して比較する工程、を備える方法。

【請求項33】 請求項30記載の方法であって、更に、  
前記検体の核酸配列とハイブリッドを形成する完全対合および不対合ブローブのハイブリッド形成強度と、基準配列とハイブリッドを形成する完全対合および不対合ブローブのハイブリッド形成強度とを、比の閾値に対して比較する工程、を備える方法。

【請求項34】 請求項30記載の方法であって、更に、  
決定行列を用いて、前記検体の核酸配列における前記遺伝子の発現変化を決定する工程、を備える方法。

【請求項35】 請求項30記載の方法であって、  
前記検体の核酸配列における前記遺伝子の発現変化が、増加、増加傾向、減少、減少傾向、変化なし、から成るグループから選択される、方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、コンピュータシステムの分野に関する。更に詳しくは、本発明は、核酸配列等の生物学的配列を解析するためのコンピュータシステムに関する。

【0002】

【従来の技術】基質上で物質の配列（アレイ）を形成し、利用するための装置やコンピュータシステムが知られている。例えば、本明細書に組み込まれるPCT出願W092/10588には、核酸等の物質の配列を決定し、あるいは、その配列をチェックする手法が述べられている。このような操作を実行するための配列（アレイ）は、本明細書に組み込まれる米国特許第5,143,854号および米国特許出願第08/249,188号に開示されているバイオニア的手法等の方法に従

て形成される。

【0003】これらにおいて説明されている手法に従って、核酸ブローブの配列（アレイ）を、基質あるいはチップ上の既知の位置に形成する。次に、蛍光標識された核酸をチップに接触させて、標識された核酸がチップに結合している位置を示すイメージファイル（イメージファイルは処理されてセルフファイルになる）をスキャナーで生成する。このセルフファイルと所定の位置におけるブローブの強度とに基づいて、DNAあるいはRNAのモノマー配列のような情報を得ることが可能になる。このようなシステムを用いることにより、例えば、脾臓癌性線維症に關係する突然変異、（所定の癌に關係する）P53遺伝子、HIV等の遺伝子特性の研究および検出に利用可能なDNAの配列（アレイ）が形成される。

【0004】

【発明が解決しようとする課題】本明細書に組み込まれる米国特許出願第08/531,137号（代理人事件番号16528X-008210）、第08/528,656号（代理人事件番号16528X-017600）、および、第08/618,834号（代理人事件番号16528-016400）に、塩基呼び出しを行うための革新的なコンピュータ利用技術が開示されている。しかし、これらのバイオニア的的手法によって現在利用され、また、利用可能になっている膨大な量の情報を評価・解析処理するためには、更に、コンピュータシステムおよび方法を改良することが必要となる。

【0005】更に、遺伝子発現を監視するためのコンピュータ利用技術の改良も必要である。多くの病態は、遺伝子DNAの複製回数の変化あるいは所定の遺伝子の（例えば、開始、RNA前駆体の供給、RNAプロセッシング等の制御による）転写レベルの変化に起因する種々の遺伝子の発現レベルにおける相違によって特徴づけることができる。例えば、遺伝物質の損失および獲得は、悪性の形態変換及び進行に重要な役割を果たす。また、（例えば、腫瘍遺伝子あるいは腫瘍抑制遺伝子等の）所定の遺伝子の発現（転写）レベルの変化は、種々の癌の存在および進行を示す指標として働く。

【0006】病気や細胞周期及び細胞発育の制御も、所定の遺伝子の転写レベルの変化によって特徴づけられる。例えば、ウイルス感染は、所定のウイルスの遺伝子の発現増強によって特徴づけられることが多い。例えば、単純ヘルペスウイルス感染、（伝染性単核症等の）エプスタイン-バーウイルス感染、サイトメガロウイルス感染、水痘-帯状ヘルペスウイルス感染、パルボウイルス感染、ヒトパピローマウイルス感染等の発生は、すべて、それぞれのウイルスに存在する種々の遺伝子の発現増強によって特徴づけられる。特徴的なウイルス遺伝子の増強発現レベルを検出することによって、その病態を効果的に診断することが可能になる。特に、単純ヘルペスウイルス等のウイルスは、長い潜伏期間を経て、

突然、短い時間で、爆発的に複製される。特徴的なウィルス遺伝子の発現レベルを検出することによって、このような活動的な増殖（おそらくは、その結果としての感染）状態を検出することが可能になる。

#### 【0007】

【課題を解決するための手段およびその作用・効果】本発明は、核酸配列等の生物学的配列（生体配列）を解析するための革新的なシステムおよび方法を提供する。コンピュータシステムは、検体配列における塩基を呼び出すために、核酸プローブと検体の核酸配列とのハイブリッド形成親和性を示すハイブリッド形成強度を解析することができる。複数の塩基コールは、単一の塩基コールを形成するために組み合わせるようにしてもよい。更に、コンピュータシステムは、遺伝子発現、あるいは、基準（基線）と比較された遺伝子発現の変化を監視するために、ハイブリッド形成強度を解析するようにしてもよい。

【0008】本発明の一態様によれば、検体の核酸配列における未知の塩基を呼び出すためのコンピュータによる方法は、複数の核酸プローブに関するハイブリッド形成強度であって、それぞれ核酸プローブと前記検体の核酸配列との間のハイブリッド形成親和性を示すハイブリッド形成強度を受け取る工程と、各組のプローブに関して前記未知の塩基に対する塩基コールを演算する工程と、前記複数のプローブに関して最も発生頻度の高い前記未知の塩基に対する塩基コールに従って、前記複数のプローブに関する単一の塩基コールを演算する工程と、を備える。典型的には、この単一塩基コールがスクリーンディスプレイ上に表示され、ユーザは、この単一塩基コールの由来である複数の塩基コールを表示するかあるいは表示しないかを選択できる。

【0009】本発明の別の態様によれば、コンピュータによる塩基呼び出し処理に関するパラメータを動的に変更する方法は、ユーザによって変更可能なパラメータを含む前記塩基呼び出し処理を利用して、検体の核酸配列の少なくとも一部に関する複数の塩基コールを生成する工程と、前記検体の核酸配列の少なくとも一部に関する前記複数の塩基コールを表示する工程と、前記塩基呼び出し処理のパラメータを表示する工程と、前記塩基呼び出し処理のパラメータに対する新しい値を特定する入力を受取る工程と、前記塩基呼び出し処理と前記パラメータに対する新しい値とを用いて、前記検体の核酸配列の少なくとも一部に対する更新された塩基コールを生成する工程と、前記検体の核酸配列の少なくとも一部に対する前記更新された塩基コールを表示する工程と、を備える。典型的には、ユーザによって変更可能なパラメータは、定数、閾値、または、範囲である。

【0010】本発明の更に別の態様によれば、検体の核酸配列において遺伝子の発現を監視するコンピュータに

よる方法は、前記遺伝子に対して完全に相補的である完全対合プローブ（パーフェクトマッチプローブ）と前記遺伝子に対して不対合な塩基を少なくとも一つ備える不対合プローブ（ミスマッチプローブ）との複数のプローブ対に関する複数のハイブリッド形成強度であって、前記完全対合および不対合プローブと前記検体の核酸配列との間のハイブリッド形成親和性を示すハイブリッド形成強度を入力する工程と、各対の完全対合プローブと不対合プローブのハイブリッド形成強度を比較する工程と、前記検体の核酸配列の遺伝子発現コールを生成する工程と、を備える。好ましい実施例においては、発現コールは、発現している、どちらともいえない、または、発現していない、として示される。

【0011】本発明の別の態様によれば、検体の核酸配列において遺伝子の発現変化を監視するコンピュータによる方法は、前記遺伝子に対して完全に相補的である完全対合プローブと前記遺伝子に対して不対合な塩基を少なくとも一つ備える不対合プローブとの複数のプローブ対に関する複数のハイブリッド形成強度であって、前記完全対合および不対合プローブと前記検体の核酸配列との間のハイブリッド形成親和性を示すハイブリッド形成強度を入力する工程と、前記検体の核酸配列の遺伝子発現レベルを生成するために、各対の完全対合プローブと不対合プローブのハイブリッド形成強度を比較する工程と、前記遺伝子発現レベルを、基準の遺伝子発現レベルと比較することにより発現変化を決定する工程と、を備える。発現の変化を、ディスプレイスクリーン上にグラフとして表示するようにしてもよい。

【0012】本発明の性質及び利点を更に理解するために、以下に、図面に基づいて、本発明を詳細に説明する。

#### 【0013】

##### 【発明の実施の形態】

##### 概要

本発明は、検体の核酸配列においてヌクレオチドを同定し（即ち、塩基を呼び出し）、遺伝子発現を監視するための革新的な方法を提供する。以下に、本発明の好適な実施例を説明する。但し、以下の実施例は、例示に過ぎず、発明の範囲を限定するものではない。

【0014】図1は、本発明のソフトウェアを実行するために用いられるコンピュータシステムの例を示す。図1に示すコンピュータシステム1は、モニタ3、スクリーン5、キャビネット7、キーボード9、および、マウス11を備える。マウス11は、マウスボタン13等の一つあるいは複数のボタンを有する。キャビネット7には、本発明のコンピュータコードを含むソフトウェアプログラムを記憶して検索するために用いられるCD-ROMドライブ15及び（図示しない）ハードドライブが収容されている。この実施例では、CD-ROM17がコンピュータ読み取り可能な媒体として示されている

10

20

30

40

50

が、フロッピーディスク、DRAM、ハードドライブ、フラッシュメモリ、テープ等、他のコンピュータ読み取り可能な媒体を用いることもできる。キャビネット7には、更に、プロセッサ、メモリ等の（図示しない）周知のコンピュータ構成部品が収容されている。

【0015】図2は、本発明のソフトウェアを実現するために用いられるコンピュータシステム1を示すシステムブロック図である。図1に示すように、コンピュータシステム1は、モニタ3とキーボード9を備える。コンピュータシステム1は、更に、中央処理装置50、システムメモリ52、入出力制御装置54、ディスプレイアダプタ56、リムーバブルディスク58、固定ディスク60、ネットワークインターフェース62、スピーカ64等のサブシステムを備える。リムーバブルディスク58は、フロッピー、テープ、CD-ROM、可動ハードドライブ、フラッシュメモリ等の取り外し可能なコンピュータ読み取り可能な媒体を示す。また、固定ディスク60は、内蔵ハードドライブ等を示す。更に、コンピュータシステム1は、本発明の実現に適した他のサブシステムを備えるようにしてもよい。例えば、コンピュータシステム1に2つ以上の処理装置50を備えるようにしてもよい（即ち、マルチプロセッサシステム）、メモリキャッシュを備えるようにしてもよい。

【0016】66等の矢印は、コンピュータシステム1のシステムバス・アーキテクチャを示す。但し、この矢印は、サブシステムを連結するあらゆる結合スキームを示すものである。例えば、ディスプレイアダプタ56がローカルバスを介して中央処理装置50に接続されるようにすることもできるし、あるいは、システムがメモリキャッシュを備えるようにすることもできる。図2に示されるコンピュータシステム1は、本発明の実現に適したコンピュータシステムの一例に過ぎず、当業者に周知のように、サブシステムの他の構成において本発明を実現することも可能である。コンピュータシステムの一例として、サン・マイクロシステムズ社のワークステーションを用いることもできる。

【0017】VLSIPS技術（商標）は、非常に小さなチップ上に、非常に大きなオリゴヌクレオチドプローブの配列（アレイ）を形成するための方法である。詳細は、本明細書に組み込まれる米国特許第5,143,854号およびPCT特許公報WO90/15070および92/10092を参照のこと。チップ上のオリゴヌクレオチドプローブを用いて、対象となる検体核酸（以下、「標的」核酸）の相補的核酸配列を検出する。

【0018】本発明は、ハイブリッド形成された核酸プローブを含むチップに関してハイブリッド形成強度ファイルを解析する方法を提供する。本実施例において、このファイルは、生物学的配列から得られる蛍光データを表しているが、放射性強度データ等の他のデータを表すようにしてもよい。即ち、本発明は、ハイブリッド形成

の蛍光測定値の解析に限定されるものではなく、ハイブリッド形成の他の測定値の解析にも容易に適用可能である。

【0019】本発明は、チップマスクを設計し、チップ上でプローブを合成し、核酸を標識して、ハイブリッド形成された核酸プローブを走査するコンピュータシステムの一部として説明する。このようなシステムは、本明細書に組み込まれる米国特許出願第08/249,188号に詳細に説明されている。但し、本発明は、例えば遠隔位置に置かれたシステムによって生成されるデータを解析するために、全体システムとは独立に利用することも可能である。

【0020】図3は、RNAやDNA等の生体物質の配列（アレイ）を形成し、解析するコンピュータ化されたシステムを示す。コンピュータ100は、RNAあるいはDNA等の生体高分子の配列（アレイ）を設計するために用いられる。ここで、コンピュータ100は、図1及び図2に示すように、例えば、適当なメモリとCPUを備え、ウィンドウズNT環境下にある適当にプログラムされたIBMパーソナルコンピュータ互換機である。コンピュータシステム100には、対象となる遺伝子の特性に関するユーザーからの情報、および、その配列（アレイ）の所望の特性に関する情報が入力される。あるいは、GenBank等の外部あるいは内部データベース102から、対象となる所定の遺伝子配列に関する情報を入力するようにしてもよい。コンピュータシステム100は、PCT出願WO92/10092に記載されているような転換行列の形で一組のチップ設計コンピュータファイル104および他の関連コンピュータファイル106を出力する。

【0021】チップ設計ファイルは、DNA等の分子配列の形成に用いられるリソグラフィックマスクを設計するためのシステム106に与えられる。システムあるいはプロセス106が、マスク110の製造に必要なハードウェア、および、効率よくマスク上にマスクパターンを配置するために必要なコンピュータハードウェア並びにソフトウェア108を備えるようにしてもよい。図3の他の特性と同様に、このような構成要素を物理的に同一の位置に配置させてもささなくてもよいが、図示の便宜上、図3では同じところに表示している。システム106は、マスク110、あるいは、クロム-ガラスマスク等、高分子配列（アレイ）の形成に用いられる他の合成パターンを生成する。

【0022】合成システム112では、マスク110およびシステム100から入力されるチップ設計に関する所定の情報が用いられる。合成システム112は、基質あるいはチップ114上に高分子配列（アレイ）を形成するために必要なハードウェアおよびソフトウェアを備える。例えば、合成装置112は、光源116および基質あるいはチップ114が設置される化学フローセル1

18を備える。マスク110を光源と基質／チップとの間に置いて、チップの所定領域から保護剤を除去するために、適当な回数、基質とマスクとの間を相対的に並進させる。所定の薬剤をフローセル118から流して、保護除去された領域への結合、および、洗浄等の操作を行う。このような操作は、全て、適当にプログラムされたコンピュータ119によって実行される構成が望ましい。ここで、コンピュータ119は、マスク設計及びマスク形成に用いられたコンピュータと同じコンピュータでも違うものでもよい。

【0023】合成システム112によって形成された基質をより小さなチップに分割して、マークされた標的に曝すようにしてもよい。ここで、標的は、基質上の一つあるいは複数の分子に対して相対的であってもよく、あるいはそうでなくてもよい。また、標的は、(図3に＊印で示す)蛍光標識等の標識でマークされ、走査システム120内に配置される。走査システム120は、適当にプログラムされたデジタルコンピュータ122によって制御される。このデジタルコンピュータ122も、合成、マスク形成、およびマスク設計に用いられたコンピュータと同じコンピュータでも違うものでもよい。スキャナ120は、標識された標的(＊)が基質に結合した位置を検出するために用いられる共焦顕微鏡あるいはCCD(電荷結合素子)等の検出装置124を備える。スキャナ120は、蛍光標識された標的の場合には、基質上の位置の開数として蛍光強度(フォトン数あるいは電圧のような他の関係する測定値)を示すイメージファイル124を出力する。標識された標的が高分子配列(アレイ)とより強固に結合している部分ではフォトン数が大きくなり、また、基質上の高分子中のモノマー配列は位置の開数として既知であるため、標的と相対的な基質上の高分子の配列を決定することができる。

【0024】イメージファイル124は、本発明の視覚化および解析方法を組み込んだ解析システム126に入力される。この解析システムも種々のコンピュータシステムのいずれでもよい。本発明は、チップ設計ファイルおよびイメージファイルを解析して、適当な出力128を与える様々な方法を提供する。また、本発明を用いて、DNAあるいはRNA等の標的における所定の突然変異を同定するようにしてもよい。

【0025】図4は、本発明の一実施例で用いられる全体的なソフトウェアシステムの概要を示す。図4に示されるように、システムは、まず、ステップ202で、所定の解析の対象となる遺伝子配列、即ち、標的を同定する。対象となる配列は、例えば、遺伝子の正常部分あるいは突然変異部分、遺伝を決定づけるあるいは修飾的情報を与える遺伝子等である。テキストファイルの手入力によって配列を選択してもよいし、あるいは、Gen Bank等の外部資源から選択するようにしてもよい。ステップ204では、システムが遺伝子の評価を行うこ

とにより、チップ上でどのプローブが望ましいかを決定し、あるいは、ユーザーが決定するための情報を与え、チップ上でどのプローブの適当な「レイアウト」を求める。

【0026】チップは、通常、既知の配列を有する標準核酸配列に対して相対的なプローブを備える。野生型プローブは、参照配列(標準配列)と理想的にハイブリッド形成するプローブであり、野生型遺伝子(チップ野生型とも称される)は、従って、チップ上の野生型プローブと理想的にハイブリッド形成する。標的配列は、突然変異、挿入、欠失等の存在を除いては、参照配列とほぼ同じである。このレイアウトは、遺伝配列の「読み込み」及び／あるいは周縁効果の最小化、合成の容易さ等を可能にするチップ上の構成を含む所望の特性を実現するものである。

【0027】図5は、チップの全体的なレイアウトを示す。チップ114は、複数のユニットから成り、各ユニットには、一つの野生型配列あるいは複数の野生型配列が様々な配置で並べられている。図にはユニット1の詳細を示してあるが、このように、各ユニットは、プローブを含むチップ上の領域である複数のセルから形成されている。各ユニットは、複数の関連するセルを含む。ここで、セルは、一つあるいは複数の分子(例えば、核酸プローブ)の多くの複製を含む基質上の領域を意味する。

【0028】各ユニットは、縦横(横は「レーン」)に配置された複数のセルから成る。例えば、5つの関連セルから成る組には、野生型セル220、「突然変異」セル222、「ブランク」セル224が含まれる。セル220は、野生型配列の一部の縮体である野生型プローブを含む。セル222は、野生型配列に対する「突然変異」プローブを含む。例えば、野生型プローブが3'-ACGTである場合、3'-ACAT、3'-ACCT、3'-ACGT、および、3'-ACTTは「突然変異」プローブである。セル224は、プローブを含まない「ブランク」セルである(「ブランク」プローブとも称される)。ブランクセルにはプローブが含まれないため、標識された標的はこの領域でチップに結合しない。このため、ブランクセルの領域を用いて、バックグラウンド強度の測定が行われる。

【0029】図4に戻って、ステップ206で、合成用のマスクが設計される。ステップ208では、ソフトウェアがマスク設計およびレイアウト情報を利用することにより、DNA等の高分子チップが形成される。このソフトウェア208は、基質とマスクとの相対的な並進、所望の薬剤のフローセル内における流れ、フローセル内の合成温度等のパラメータを制御する。ステップ210では、他のソフトウェアを用いて、以上のようにして合成され、標識された標的に曝されたチップを走査する。ソフトウェアは、チップの走査を制御して、後に配列情

報を得るために用いられるファイルに得られたデータを格納する。

【0030】ステップ212で、コンピュータシステムは、レイアウト情報と蛍光情報とを用いて、チップ上でハイブリッド形成された核酸プローブの評価を行う。DNAチップから得られる重要な情報としては、突然変異標的の同定、及び、特定の標的の遺伝子配列の決定が挙げられる。

【0031】図6は、特定の標的DNAがDNAプローブ114の配列（アレイ）に結合している状態を示す。この簡単な例に示されるように、配列中に以下のプローブが形成される（野生型プローブに関しては一つのプローブだけが示されている）。

【0032】

```
3'-AGAACGT
AGACCGT
AGAGCGT
AGATCGT
...
```

【0033】これらのプローブでは、それぞれ、一つの塩基のみが異なっている。即ち、調査位置において単一の塩基不対合が存在する。従って、核酸配列中のこの調査位置における塩基を同定するように、プローブが設計される。即ち、ここでは、1つのユニットとは、複数組の関連するプローブのことを言い、各組には、調査位置における単一の塩基不対合が異なる複数のプローブが含まれる。

【0034】5'-TCTTGCAの配列を有する蛍光標識された（あるいは別の方法でマークされた）標的をDNAプローブの配列（アレイ）に曝した場合、これは、3'-AGAACGTプローブのみに相補的であり、フルオレセインは、3'-AGAACGTが位置するチップ表面上で主に観察される。一つの塩基のみがそれぞれ異なる各組のプローブに関するイメージファイルには、各々プローブに対応する4つの蛍光強度が含まれる。即ち、蛍光強度は、各々、他のプローブとは異なる各プローブのヌクレオチドあるいは塩基に結びつけられる。更に、イメージファイルには、バックグラウンドの蛍光強度として用いられる「ブランク」セルが含まれる。特定の塩基位置に関する5つの蛍光強度を解析することにより、ここに開示される本発明の方法を用いて、このような配列（アレイ）から配列（シーケンス）情報を得ることができる。

【0035】図7は、チップ上に列で配置されたプローブを示す。参照配列（即ち野生型チップの配列）を下つきの数字で示された5つの調査位置と共に示す。調査位置は、多くの場合、標的配列が突然変異を含む、即ち、標的配列が参照配列と異なっているような、参照配列中の塩基位置である。チップには、それぞれの調査位置に対応する5つのプローブセルが含まれる。プローブセル

には、各々、調査位置に共通の塩基を有するプローブの組が含まれる。例えば、参照配列は、第一の調査位置1、に塩基Tを有する。この調査位置に関する野生型のプローブは、3'-TGACであり、プローブ中の塩基Aは、参照配列における調査位置の塩基に対して相補的である。

【0036】この第一の調査位置1、に関しては、4つの「突然変異」プローブセルが存在する。4つの突然変異プローブは、それぞれ、3'-TGAC、3'-TGCC、3'-TGGC、3'-TGTGである。4つの突然変異プローブは、それぞれ調査位置において一つの塩基のみが異なっている。図示するように、野生型のプローブと突然変異プローブはチップ上に列配置されている。突然変異プローブのうちの一つ（この場合には、3'-TGAC）は、野生型のプローブと等しいため、突然変異を証明するものとはならない。但し、図8に示すように、冗長性により突然変異を視覚的に示すことができる。

【0037】図7に示すように、チップは、他の調査位置1、-4、の各々に関しても、野生型プローブと突然変異プローブを有している。それぞれの場合において、野生型プローブは突然変異プローブの一つと等価である。

【0038】図8は、チップ上の標的と図7の参照配列とのハイブリッド形成パターンを示している。比較のために、チップの一番上に参照配列を示している。チップには、WT列（野生型）、A列、C列、G列、T（またはU）列が含まれる。各列は、プローブを含むセル列である。WT列のセルは、参照配列に相補的なプローブを含む。A列、C列、G列、T列のセルは、それぞれ、列の名称となっている塩基が調査位置に存在することを除けば参照配列に対して相補的であるプローブを含んでいる。

【0039】一実施例においては、セル内のプローブのハイブリッド形成は、標識された標的配列の結合に起因するセルの蛍光強度（例えば、光子数）によって決定される。蛍光強度はセルによって大きく異なる。単純化するために、図8では、領域を塗りつぶすことによって、セルが高度のハイブリッド形成をしていることを示す。WT列を観察すると、調査位置1、において野生型セルの領域が塗りつぶされていないため、この位置1、に突然変異が存在することが容易にわかる。C列のセルは塗りつぶされているため、TからGへの突然変異であることがわかる（突然変異プローブは相補的であるため、CセルはG突然変異を示す）。好ましい実施例では、WT列は利用されず、調査位置の塩基を呼び出すために、「ブランク」セルを含まない4つのセルが使用される。

【0040】実際には、突然変異が存在する調査位置近傍のセルの蛍光強度は相対的に低く、突然変異の周囲は

「暗い領域」となる。このように蛍光強度が低くなるのは、突然変異近辺の調査位置のセルは、標的配列と完全に相補的なプローブを含んでおらず、その結果、これらのプローブと標的配列とのハイブリッド形成率が低いからである。例えば、調査位置I<sub>1</sub>及びI<sub>2</sub>において、これらの位置に存在するプローブはどれも標的配列に相補的でないため、セルの相対的強度が比較的低くなる。蛍光強度が低いとデータの解像度も低くなるが、本発明の方法を用いることによって、突然変異周囲の暗い領域でも精度よく塩基を呼び出すことができ、暗い領域内に存在する他の突然変異を同定することができる。

【0041】図9は、チップ上における標準タイル構造及び代替タイル構造を示す。図示するように、チップには12のユニット（ユニット<sub>1,12</sub>）が含まれている。ユニット<sub>1,12</sub>は、同一の参照配列に相補的なプローブを含むようにタイル状に構成されている（即ち、チップ上で設計され、合成されている）。以下、このユニットグループを標準グループと称する。本発明では、他に明記しない限り、標的配列に関する塩基コールには標準グループが用いられる。

【0042】ユニット<sub>1,12</sub>は、標準グループの参照配列とは異なった同一の参照配列に相補的であるプローブを含むようにタイル状に構成される。以下、このユニットグループを代替グループと称する。ユニット<sub>1,12</sub>は、標準グループ及び第一の代替グループの参照配列とは異なった参照配列に基づく別の代替グループを構成している。これらの参照配列は、互いに異なっているが、多くの場合、かなり似ている。例えば、参照配列を、少しずつ異なったHIVの突然変異としてもよい。本発明の実施例は、標的配列の塩基コールには通常用いられないような参照配列に基づくタイル構成からの情報を評価し、利用する。

【0043】あるグループ内の複数のユニットは同一のプローブを備えるようにしてもよいし、異なった構造のプローブを備えるようにしてもよいし、あるいは、一つまたは複数の別のチップに由来するプローブを備えるようにしてもよい。例えば、1つのユニットが、プローブの第三位置に調査位置を有する5量体プローブを備えるようにしてもよい。また、別のユニットが、第六位置に調査位置を有する10量体プローブを備えるようにしてもよい。更に、これらのユニットを同一のチップ上に配置してもよいし、異なったチップ上に配置してもよい。

【0044】図9左下の拡大図に示すように、ユニットの各ブロックには、通常、A、C、G、Tで示される4つのセルが含まれる。セル内の各プローブの調査位置にどの塩基があるかを、この塩基表示によって示している。各セルには、普通、何百、何千という同一の核酸プローブが存在する。

【0045】好適な実施例では、セルは、参照配列に沿って、順番に互いに隣接して配置されているが、チップ

上での位置さえ確定されている限り、特定の位置にセルを配置しなければならないという決まりはない。更に、実験を首尾一貫したものにするためには、単一のチップ上で異なったグループを合成することが望ましいが、本発明の方法を、異なったチップ上の異なった構成から得られるデータに適用することもできる。

#### 【0046】標的配列の解析

図10は、チップから得られたハイブリッド形成強度を表示する画面である。解析が行われる場合、システムには、ハイブリッド形成されたチップを走査した画像を含むイメージファイルが入力される。ここで、イメージファイルが、蛍光強度と、標識された標的核酸配列あるいはその切片がチップに結合する位置を示すようにする構成も望ましい。

【0047】画面表示（スクリーンディスプレイ）260には、周知のウィンドウ用GUI（グラフィカルユーザーインターフェース）が用いられている。ユーザーの選択に従って、イメージファイルが表示される。ユーザーがイメージファイルの表示を選択すると、イメージファイルを含むウィンドウ262が表示される。図示されているイメージファイルには、A列、C列、G列、T列の複数の列が含まれる。

【0048】ユーザーが表示されているイメージファイル上にカーソルを動かすと、ステータスバー264に、カーソルのX、Y座標の位置とその位置における蛍光強度が示される。また、ポインティング・デバイスを用いて、イメージファイルの矩形の領域を選択することにより、サブイメージを処理することも可能である。例えば、サブイメージを拡大して、個々のセルがもっとはっきりと見えるようにすることもできる。更に、強度のコントラストを調整することによって、現在のコントラストの設定でははっきりしないハイブリッド形成強度の差異を明確にすることもできる。

【0049】図11は、関連するプローブのハイブリッド形成強度に基づいて塩基呼び出し処理を行う方法を示すフローチャートである。ここで、「関連するプローブ」とは、調査位置におけるヌクレオチド塩基が異なるプローブを示す。通常、これらのプローブは調査位置を除いて同一の構成をしているが、更に、別の塩基位置でも異なる構成にしてもよい。即ち、関連するプローブでは、少なくとも一つの塩基が異なっている。

【0050】ステップ302で、バックグラウンド強度、即ち「ブランク」セル強度を差し引くことによって、4つの関連するプローブのハイブリッド形成強度を調整する。引き算の結果ハイブリッド形成強度がゼロ以下になる場合には、小さな正の数をハイブリッド形成強度に設定することが望ましい。このように設定することにより、以降の計算で、ゼロあるいは負の数での割り算を避けることができる。

【0051】ステップ304では、ハイブリッド形成強

10

20

30

40

50



度を強度順に並べる。次に、ステップ306で、最大強度を、所定のバックグラウンド差分限界値と比較する。バックグラウンド差分限界値は、未知の塩基を正しく呼び出すために、最高強度を有するプローブが超えなければならないハイブリッド形成強度を示す数値である。即ち、バックグラウンド調整後の塩基強度は、バックグラウンド差分限界値より大きくなければならない。そうでない場合には、未知の塩基を正確に呼び出すことができない可能性がある。

【0052】関連するプローブで最も高いハイブリッド形成強度がバックグラウンド差分限界値以下の場合には、ステップ308で、(不十分な強度を示す)コードNを未知の塩基に割り当てる。そうでない場合には、ステップ310で、最高のハイブリッド形成強度と二番目に高いハイブリッド形成強度との比を計算する。

【0053】次に、ステップ312で、ステップ310で計算した比を、所定の比限界値と比較する。比限界値は、未知の塩基を同定するために必要な比を示す数値である。好ましい実施例においては、比限界値は1.2である。計算された比が比限界値より大きければ、最高のハイブリッド形成強度を有するプローブに従って、未知の塩基が呼び出される。具体的には、ステップ314で、最高強度のプローブ内の調査位置における塩基の補体を未知の塩基として呼び出す。一方、計算された比が比限界値以下である場合には、ステップ316で、二番目に高いハイブリッド形成強度と三番目に高いハイブリッド形成強度との比を計算する。

【0054】次に、ステップ318で、ステップ316で計算した比を比限界値と比較する。計算された比が比限界値より大きい場合には、ステップ320で、一番目と二番目に高いハイブリッド形成強度を有するプローブの調査位置における塩基の補体を規定する不確定コードを、未知の塩基として呼び出す。一方、計算された比が比限界値以下である場合には、ステップ322で、三番目に高いハイブリッド形成強度と四番目に高いハイブリッド形成強度との比を計算する。

【0055】次に、ステップ324で、ステップ322で計算した比を比限界値と比較する。計算された比が比限界値より大きい場合には、ステップ326で、一番目ないし三番目に高いハイブリッド形成強度を有するプローブの調査位置における塩基の補体を規定する不確定コードを、未知の塩基として呼び出す。一方、計算された比が比限界値以下である場合には、ステップ328で、(不十分な区別を示す)コードXを未知の塩基に割り当てる。

【0056】図12は、関連するプローブのハイブリッド形成強度に基づいて塩基呼び出し処理を行う他の方法を示すフローチャートである。このフローチャートでは、関連するプローブが示すハイブリッド形成強度に基づいて処理が行われる。即ち、塩基コール処理により、

調査位置における一つの塩基の不对合だけが互いに異なる複数のプローブ内の調査位置に対応するような、標的内の塩基が呼び出される。ステップ402で、システムは、標的配列に対して最高のハイブリッド形成強度を有するプローブが一つだけ存在するかどうかを判定する。存在しない場合には、不確定を意味するNを塩基に割り当てる。例えば、2つのプローブが同じ最高強度を有している場合(即ち、同強度である場合)、その塩基にNを割り当てる。

【0057】標的に対して最高のハイブリッド形成強度を有するプローブが一つだけ存在する場合には、ステップ406で、そのプローブに従って塩基の呼び出しが行われる。プローブは標的配列に相補的であるため、プローブの調査位置における塩基に相補的な塩基(C/G、A/T)を呼び出すようにしてもよい。

【0058】次に、ステップ408で、システムは、塩基コール(呼び出された塩基)が突然変異であるかどうかを判定する。即ち、参照配列の塩基と異なっているかどうかを判定する。塩基コールが突然変異塩基コールでない場合には、塩基コールを確定する。一方、塩基コールが突然変異塩基コールの場合には、システムは、ステップ410で所定の「突然変異」条件を満たしているかどうかを判定し、満たしていない場合には、ステップ412で、Nを塩基に割り当てる。

【0059】以下に突然変異条件の例を説明するが、説明をわかりやすくするために、関連するプローブのハイブリッド形成強度を次のようにラベルする。HighIntは最高のハイブリッド形成強度を、SecondIntは二番目に高いハイブリッド形成強度を、ThirdIntは三番目に高いハイブリッド形成強度を、LowIntは最も低いハイブリッド形成強度をそれぞれ示す。

【0060】例えば、呼び出された塩基が突然変異であると規定する突然変異条件は、以下の3つのテスト全てを満足させるものでなければならない。第一のテストは、HighIntとSecondIntとの差が差分限界値より大きいかどうかを判定するものである。即ち、システムは、 $(HighInt - SecondInt)$ が所定の値より大きいかどうかを判定する。この所定の値は、最高のハイブリッド形成強度が次に高いハイブリッド形成強度よりも所望の量だけ大きい場合にのみ、呼び出された塩基が突然変異塩基の可能性があると認めるように選択されたものである。

【0061】第二のテストは、第一の比が第一の比限界値未満であるかどうかを判定するものである。第一の比は以下のように規定される。

【0062】 $\{SecondInt - \sqrt{(ThirdInt * LowInt)}\} / \{HighInt - \sqrt{(ThirdInt * LowInt)}\}$

【0063】システムは、この第一の比が所定の値未満であるかどうかを判定する。この所定の値は、2つのより低いハイブリッド形成強度を差し引いた後でも、最高のハイブリッド形成強度が、次に高いハイブリッド形成



強度に対して所望の比より大きい場合のみ、呼び出された塩基が突然変異塩基の可能性があると認めるように選択されたものである。

【0064】第三のテストは、隣接比が隣接比限界値より大きいかどうかを判定するものである。隣接比は以下のように規定される。

【0065】 $\text{HighInt}_n / \{\text{HighInt}_n - \sqrt{(\text{HighInt}_{n-1} * \text{HighInt}_{n+1})}\}$

【0066】ここで、添え字のnは、呼び出されている対象の塩基位置における値を示す。また、n+1及びn-1は、隣接する塩基位置における値を示す。システムは、この隣接比が所定の値より大きいかどうかを判定する。この所定の値は、最高のハイブリッド形成強度が、隣接する最高のハイブリッド形成強度を引いた後の最高ハイブリッド形成強度に対して所望の比より大きい場合のみ、呼び出された塩基が突然変異塩基の可能性があると認めるように選択されたものである。

【0067】このように、突然変異条件の全てが満たされた場合のみ、呼び出された塩基が突然変異塩基と規定される。突然変異はかなりまれであるため、突然変異が存在する可能性が高い場合にのみ突然変異塩基を呼び出すようにする。突然変異条件が満たされない場合には、呼び出された塩基は不確定であるか、あるいは、標準塩基と同一である（統計的に、正しい塩基コールの可能性はある）。

【0068】この実施例では、3つの突然変異条件を用いたが、突然変異条件の一つのみ（例えば、上述の条件のどれか一つ）とすることもできる。また、参照文献として先に挙げた米国特許出願に記載されているような塩基呼び出し方法を利用することもできる。

【0069】図13は、一つのユニットグループに関して塩基呼び出し処理を行う方法を示すフローチャートである。上述したように、一つのユニットには、複数組の関連セルが含まれ、関連セルは、調査位置における一つの塩基がそれぞれ異なっているプローブを含む。システムは、まず、（例えば、ハイブリッド形成されたチップを走査するスキャナによって生成されたイメージデータファイルから）ハイブリッド形成強度を入力し、更に、このハイブリッド形成強度に対応するプローブの構造を入力する。次に、バックグラウンド強度（例えば、「ブランク」セルの強度、あるいは、プローブを含まないチップの別の領域において測定された強度）を入力されたハイブリッド形成強度から引く。ここで、バックグラウンド引き算後のハイブリッド形成強度の最低値を1（フォトン数1）に設定しておいてもよい。

【0070】ハイブリッド形成強度は、プローブ（あるいはプローブの複数組の複製）と標的配列との間で測定されるハイブリッド形成の度合いを意味する。例えば、測定されたセルのフォトン数の平均をハイブリッド形成強度としてもよい。ここで、フォトン数は、セル内のプロ

ープに結合し、フルオレセインで標識された標的配列に由来する。

【0071】ステップ452で、システムは、コール処理の対象となる標的配列の塩基位置を入力する。次に、ステップ454で、グループの各ユニットに関してその塩基位置における塩基コール処理を実行する。具体的には、その塩基位置における各ユニットの関連セルのハイブリッド形成強度を解析する。この解析を行うことにより（解析処理の詳細に関しては、図11及び図12で既に説明した）、システムは、各ユニットに関する塩基コールを求める。例えば、グループ内に5つのユニットが存在する場合には、5つの塩基コールが生成される。

【0072】システムは、ステップ456で、そのグループのユニットに関する複数の塩基コールを解析して、そのグループの塩基コールを一つ決定する。具体的には、そのグループのユニットから呼び出される頻度が一番高かった塩基を、そのグループの塩基として呼び出すように構成してもよい。例えば、5つのユニットが存在し、各ユニットが次のような塩基コールがなされたとする。

【0073】「T」-3ユニット

「G」-1ユニット

「N」-1ユニット

【0074】この場合には、5ユニットのうち3ユニットが呼び出しているTが、塩基として呼び出される。同じ頻度で呼び出された塩基が複数ある場合には、塩基を呼び出すユニットの最高平均ハイブリッド形成強度等の他のファクターを考慮する。本発明の実施例においては、図15に示される方法が用いられる。

【0075】ステップ458で、解析を実行すべき次の塩基位置が存在するかどうかを判定する。本発明の方法を適用して、標的核酸配列の全ての塩基位置に関して塩基コール処理を実行することもできるし、不必要な塩基位置をスキップして、所定の塩基位置においてのみ塩基コール処理を実行することもできる。

【0076】図14は、複数のユニットグループに関して塩基呼び出し処理を行う方法を示すフローチャートである。図9に示したように、解析対象となる一つあるいは複数のチップ上に複数のグループが存在する場合がある。複数のグループは、異なった参照配列に従って配置されていてもよい。但し、これは、複数のグループのハイブリッド形成に関する情報が利用されない可能性を示すものではない。通常、標準グループの参照配列が標的配列と一致する可能性が最も高いが、代替グループの一つの方が一致の可能性が高い場合には（即ち、塩基コール処理により適している場合には）、そのグループを用いて塩基処理を行う。

【0077】ステップ502で、システムは、標準グループ及び代替グループのユニットに関して塩基コール処理を実行する。塩基コール処理は、例えば、先に図13

で説明したように行われる。

【0078】次に、ステップ504で、各ユニットグループに関して塩基コールを一つ求める。ここで、ユニットからの呼び出し頻度が最も高かった塩基を塩基コールとしてもよいし、あるいは、図15に従って後に詳細を説明する方法で、各グループの塩基コールを決定するようにしてもよい。

【0079】各ユニットグループ（即ち、標準グループと代替グループ）に関して一つの塩基コールをもとめた後、ステップ506で、塩基位置が入力される。次に、入力された塩基位置に関して塩基コールを実行するために最適なユニットグループを選択する。具体的には、例えば、調査位置近傍、即ち、調査位置周囲のウィンドウにおいて、標的配列に対する不適合が最も少ないグループの参照配列を求めることにより、最適グループの選択が行われる。調査位置近傍において不適合が最も少ないユニットグループは、最も精度の高い塩基コールを生成する可能性が高い。最適グループを選択する処理に関する詳細は、図16を参照して後述する。

【0080】次に、ステップ510で、最適のユニットグループに従って、入力された塩基位置における塩基の呼び出し処理を行う（即ち、ステップ504で決定されたそのグループの最適コールを用いる）。塩基コールが決定されると、処理は次のステップに移り、塩基コール処理を実行すべき次の塩基位置が存在するかどうかの判定が行われる。次の塩基位置が存在する場合には、ステップ506に戻り、次の塩基位置に関して塩基コール処理を実行する。

【0081】図15は、一つのユニットグループに関して塩基呼び出し処理を行う方法を示すフローチャートである。ステップ602で、システムは、所定の塩基位置において大部分のユニットが同じ塩基を呼び出しているかどうかを判定する。この場合、（不確定を示すコードNが割り当てられたものを除き）塩基を呼び出すユニットのみが判定の対象となる。例えば、7つのユニットが存在し、各ユニットが次のような塩基コールを実行したとする。

【0082】「G」-3ユニット

「T」-1ユニット

「N」-4ユニット

【0083】この場合には、4つの確定塩基コールのうち3つがGを呼び出しているため、まず、そのユニットグループに関してGを塩基として呼び出す。例外規則が適用される場合を除いて、ステップ604で、この塩基が多数塩基として呼び出される。

【0084】例外規則は、そのユニットグループに関してどの塩基を呼び出すかを定める条件を規定したものである。この例外規則は、多数塩基コールを変える条件や、大多数のユニットによって呼び出される塩基が存在しないような状況に対処する条件を規定する。例外規則

の一例として、（一つのユニットがある塩基を呼び出し、他のユニットが別の塩基を呼び出すような場合に適用される）隣接するプローブのハイブリッド形成強度を解析するタイブレーク規則が挙げられる。また、別の例として、3つのユニットがそれぞれ別の塩基を呼び出し、そのうちの一つが標準塩基を呼び出すものである場合、標準塩基をそのユニットグループの塩基コールとする例外規則が挙げられる。他の例外規則に関しては、前述に示す。

10 【0085】ステップ606で、例外規則が適用されるかどうかの判定がなされる。例外規則が適用されると判定された場合には、ステップ608で例外規則が適用される。

【0086】図16は、塩基コール処理を実行するために最適のユニットグループを選択する方法を示すフローチャートである。上述したように、調査位置近傍において、標的配列に対する不適合が最も少ないグループの参照配列を求めることにより、最適グループの選択が行われる。調査位置近傍において不適合が最も少ないユニットグループは、最も精度の高い塩基コールを生成する可能性が高い。解析の対象となる調査位置近傍のウィンドウは、所定の値に設定してもよく、あるいは、プローブ構造に従って設定してもよい。例えば、全てのグループのプローブに関して、調査位置からの最大距離が、調査位置の一方の側に対しては8塩基位置、また、もう一方の側に対しては10塩基位置である場合、この範囲の塩基位置を含むようにウィンドウを設定してもよい。

【0087】ステップ702で、システムは、標準ユニットグループおよび代替ユニットグループに関して不適合得点を計算する。不適合得点は、標的配列に対する参照配列の不適合の数を示すものである。不適合得点を求めるために、参照配列のうち少なくとも2つが異なっている塩基位置のみを解析するようにしてもよい。全ての参照配列が或る塩基位置において同一である場合には、この塩基配列をスキップすることができる。

【0088】少なくとも2つの参照配列が異なっている各塩基位置において、あるグループに関して呼び出された塩基（標的配列における有望な塩基を示す塩基コール）が参照配列の対応する塩基と異なっているかどうかを判定する。ここで、塩基コールと参照配列の対応する塩基とが異なっている場合には、不適合得点に1が加えられる。各グループに関する不適合得点は、最初にゼロに設定される。

【0089】不適合得点は、標的配列と異なっている参照配列部分の塩基位置の数を示すものだと考えることができる（全ての参照配列において同一である塩基位置を除くようにしてもよい）。この概念をさらに説明するために、以下に簡単な例を挙げる。この例では、以下のようないつの標準グループと2つの代替グループが存在する。

で説明したように行われる。

【0078】次に、ステップ504で、各ユニットグループに関して塩基コールを一つ求める。ここで、ユニットからの呼び出し頻度が最も高かった塩基を塩基コールとしてもよい。あるいは、図15に従って後に詳細を説明する方法で、各グループの塩基コールを決定するようにしてもよい。

【0079】各ユニットグループ（即ち、標準グループと代替グループ）に関して一つの塩基コールをもとめた後、ステップ506で、塩基位置が入力される。次に、入力された塩基位置に関して塩基コールを実行するために最適のユニットグループを選択する。具体的には、例えば、調査位置近傍、即ち、調査位置周囲のウィンドウにおいて、標的配列に対する不適合が最も少ないグループの参照配列を求めることにより、最適グループの選択が行われる。調査位置近傍において不適合が最も少ないユニットグループは、最も精度の高い塩基コールを生成する可能性が高い。最適グループを選択する処理に関する詳細は、図16を参照して後述する。

【0080】次に、ステップ510で、最適のユニットグループに従って、入力された塩基位置における塩基の呼び出し処理を行う（即ち、ステップ504で決定されたそのグループの最適コールを用いる）。塩基コールが決定されると、処理は次のステップに移り、塩基コール処理を実行すべき次の塩基位置が存在するかどうかの判定が行われる。次の塩基位置が存在する場合には、ステップ506に戻り、次の塩基位置に関して塩基コール処理を実行する。

【0081】図15は、一つのユニットグループに関して塩基呼び出し処理を行う方法を示すフローチャートである。ステップ602で、システムは、所定の塩基位置において大部分のユニットが同じ塩基を呼び出しているかどうかを判定する。この場合、（不確定を示すコードNが割り当てられたものを除き）塩基を呼び出すユニットのみが判定の対象となる。例えば、7つのユニットが存在し、各ユニットが次のような塩基コールを実行したとする。

【0082】「G」-3ユニット

「T」-1ユニット

「N」-4ユニット

【0083】この場合には、4つの確定塩基コールのうち3つがGを呼び出しているため、まず、そのユニットグループに関してGを塩基として呼び出す。例外規則が適用される場合を除いて、ステップ604で、この塩基が多数塩基として呼び出される。

【0084】例外規則は、そのユニットグループに関してどの塩基を呼び出すかを定める条件を規定したものである。この例外規則は、多数塩基コールを変える条件や、大多数のユニットによって呼び出される塩基が存在しないような状況に対処する条件を規定する。例外規則

の一例として、（一つのユニットがある塩基を呼び出し、他のユニットが別の塩基を呼び出すような場合に適用される）隣接するプローブのハイブリッド形成強度を解析するタイブレーク規則が挙げられる。また、別の例として、3つのユニットがそれぞれ別の塩基を呼び出し、そのうちの一つが標準塩基を呼び出すものである場合、標準塩基をそのユニットグループの塩基コールとする例外規則が挙げられる。他の例外規則に関しては、前述に示す。

【0085】ステップ606で、例外規則が適用されるかどうかの判定がなされる。例外規則が適用されると判定された場合には、ステップ608で例外規則が適用される。

【0086】図16は、塩基コール処理を実行するために最適のユニットグループを選択する方法を示すフローチャートである。上述したように、調査位置近傍において、標的配列に対する不適合が最も少ないグループの参照配列を求めることにより、最適グループの選択が行われる。調査位置近傍において不適合が最も少ないユニットグループは、最も精度の高い塩基コールを生成する可能性が高い。解析の対象となる調査位置近傍のウィンドウは、所定の値に設定してもよく、あるいは、プローブ構造に従って設定してもよい。例えば、全てのグループのプローブに関して、調査位置からの最大距離が、調査位置の一方の側に対しては8塩基位置、また、もう一方の側に対しては10塩基位置である場合、この範囲の塩基位置を含むようにウィンドウを設定してもよい。

【0087】ステップ702で、システムは、標準ユニットグループおよび代替ユニットグループに関して不適合得点を計算する。不適合得点は、標的配列に対する参照配列の不適合の数を示すものである。不適合得点を求めるために、参照配列のうち少なくとも2つが異なっている塩基位置のみを解析するようにしてもよい。全ての参照配列が或る塩基位置において同一である場合には、この塩基配列をスキップすることができる。

【0088】少なくとも2つの参照配列が異なっている各塩基位置において、あるグループに関して呼び出された塩基（標的配列における有望な塩基を示す塩基コール）が参照配列の対応する塩基と異なっているかどうかを判定する。ここで、塩基コールと参照配列の対応する塩基とが異なっている場合には、不適合得点に1が加えられる。各グループに関する不適合得点は、最初にゼロに設定される。

【0089】不適合得点は、標的配列と異なっている参照配列部分の塩基位置の数を示すものだと考えることができる（全ての参照配列において同一である塩基位置を除くようにしてもよい）。この概念をさらに説明するために、以下に簡単な例を挙げる。この例では、以下のようないつの標準グループと2つの代替グループが存在する。

10

20

30

40

50

で、簡単に解析を行うことができるように、参照配列と検体配列とが自動的に並べられる。

【0101】図17は、スクリーンディスプレイの配置を簡単に示すものであるが、図18に示すように、チップファイル及び複台ファイルから入力された情報を出さないようにしたり（非表示にしたり）、まとめたりすることもできる。例えば、ユーザーが、スクリーン領域808において、複台ファイルのファイル名の前につけられているスクリーンアイコンの+マークを「クリックする」、即ち起動した場合には、その複台ファイルに関する詳細な情報が表示される。図示されるように、チップファイルから入力された情報を組み合わせるために用いられた方法を、個々のチップファイルと共に表示するようにしてもよい。

【0102】また、ユーザーが、スクリーン領域808において、チップファイルのファイル名の前につけられているスクリーンアイコンの+マークを起動した場合には、塩基の呼び出しに用いられた工程あるいは方法を含むチップファイルに関する詳細な情報が表示される。図18の例では、図10に基づいて説明した「比に基づくアルゴリズム」が塩基呼び出し処理に用いられている。さらに、ユーザーは、スクリーン上に表示されている塩基コールに直ちに反映される塩基呼び出し処理パラメータを変更することもできる。例えば、この例では、比の限界値（「Ratio」）が1.2に設定されているが、ユーザーがこの比の限界値を1.4に増加させた場合には、そのチップに関する塩基コールを再演算して、新しい塩基コールをスクリーン領域812に表示する。塩基コール処理パラメータは、定数、閾値、範囲等のいずれの値でもよい。

【0103】図18に示すように、ユーザーが理解しやういように、（様々な構成を含む）複数の実験から得られたデータを組み合わせることができる。図17に示される検体配列440-2Aは、図18では拡大されて、この塩基コールが複数の実験から得られたものであることがわかる。ここで、複数の実験から得られたデータは、一つのチップに関するものでも、あるいは、複数のチップに関するものでもよい。言いかえると、図17および図18で検体配列440-2Aとして示されるヌクレオチド配列は、一つの実験結果を示すものではなく、複数の実験結果を組み合わせ、あるいは、まとめたものである。図18に示すように、ユーザーは、それぞれの実験から得られたデータをまとめることができる。図18の例では、それぞれの塩基に関するハイブリッド形成強度が表示されている。塩基位置100に示すように、解析中の塩基位置を反転表示することもできる。

【0104】図17に示すように、検体配列440-2Aの名称の前にもスクリーンアイコンの+マークが表示されている。このスクリーンアイコンを起動することにより、複台塩基コールを構成する個々の塩基コールを表

示することができる。また、図18に示すように、複台塩基コールは、複数の塩基コールから誘導されている。この複数の塩基コールは、塩基位置に従って、複台塩基コールと並べられる。本発明の構成により、ユーザーは、配列を解析する際に、必要に応じて、データを表示したり、非表示にしたり、あるいは、まとめたりすることができる。

#### 【0105】遺伝子発現の監視

図19は、複数対の完全対合プローブと不対合プローブのハイブリッド形成強度を比較することにより遺伝子の発現を監視する方法を示すフローチャートである。「完全対合プローブ」は、所定の標的配列に完全に相補的な配列を有するプローブを意味する。テストプローブは、通常、標的配列の一部（サブ配列）に完全に相補的である。また、「不対合対照」あるいは「不対合プローブ」は、所定の標的配列に完全に相補的でないように故意に選択された配列を有するプローブを意味する。高密度配列（アレイ）における不対合（MM）対照には、それぞれ、同じ所定の標的配列に対して完全に相補的な対応する完全対合（PM）プローブが存在する。

【0106】基質表面あるいはチップ表面に望ましくは共有結合的に結合されている、複数対の完全対合プローブ及び不対合プローブのハイブリッド形成強度が比較される。核酸プローブの密度としては、基質1cm<sup>2</sup>当たり約60個以上の異なる核酸プローブが含まれる構成が最も望ましい。図19のフローチャートでは、各ステップが順に示されているが、必ずしも、これらのステップをこの順番で実行しなくてもよい。当業者に周知のように、本発明の範囲から逸脱することなく、これらのステップの順番を変えたり、組み合わせたり、あるいはいくつかのステップを除くこともできる。

【0107】まず最初に、標的配列（あるいは遺伝子）に相補的な核酸プローブを選択する。これらのプローブは完全対合プローブである。次に、標的配列に対して完全に相補的でないように構成されたプローブを選択する。これらのプローブは不対合プローブであり、不対合プローブには、それぞれ、完全対合プローブに対して少なくとも一つのヌクレオチド不対合が含まれる。従って、不対合プローブおよびそれに対応する完全対合プローブにより、一対のプローブが形成される。上述したように、不対合プローブの中心付近にヌクレオチド不対合が存在するような構成が望ましい。

【0108】完全対合プローブの長さは、通常、標的配列に対する高いハイブリッド形成親和性を示すように選択される。例えば、核酸配列を全て20置体としてもよい。また、不確定さを排除する等様々な理由により、基質上で、異なる長さのプローブを合成するようにしてもよい。

【0109】上述したように、標的配列は、分割され、標識されて、核酸プローブを含む基質に曝される。核酸

プローブのハイブリッド形成強度を測定して、コンピュータシステムに入力する。ここで用いられるコンピュータシステムは、基質のハイブリッド形成を指示するシステムと同じものでも異なったものでもよい。いずれのコンピュータシステムを用いるにしても、遺伝子名、遺伝子配列、プローブ配列、基質上でのプローブ位置等、他の実験の詳細もシステムに入力される。

【0110】図19に示すように、ハイブリッド形成後、ステップ902で、コンピュータシステムに、複数対の完全対合プローブと不対合プローブのハイブリッド形成強度が入力される。ハイブリッド形成強度は、核酸プローブと（遺伝子に対応する）標的核酸との間のハイブリッド形成親和性を示すものである。各対には、標的核酸の一部に完全に相補的な完全対合プローブと、少なくとも一つのヌクレオチドが完全対合プローブと異なっている不対合プローブとが含まれる。

【0111】ステップ904で、コンピュータシステムが、各対の完全対合プローブと不対合プローブのハイブリッド形成強度を比較する。遺伝子が発現している場合には、完全対合プローブのハイブリッド形成強度（即ち親和性）は、対応する不対合プローブの強度よりも有意に高い。一般に、対のプローブのハイブリッド形成強度がほとんど同じであれば、遺伝子が発現されていないと考えられる。但し、一対のプローブに基づいて判定を行うわけではなく、多数対のプローブの解析に基づいて、遺伝子が発現しているかどうかの判定が行われる。複数対のプローブのハイブリッド形成強度を比較する具体的な方法の詳細に関しては、図20を参照して説明する。

【0112】完全対合プローブと不対合プローブのハイブリッド形成強度を比較した後、ステップ906で、遺伝子の発現状態を表示する。例えば、遺伝子が存在している（発現している）、どちらともいえない、あるいは、存在していない（発現していない）ことを示す発現コールをユーザーに表示するようにしてもよい。

【0113】図20は、決定行列を用いて遺伝子が発現しているかどうかを判定する方法を示すフローチャートである。ステップ952で、コンピュータシステムに、N対の完全対合プローブと不対合プローブに関する生の走査データを入力する。好適な実施例において、ハイブリッド形成強度は、基質上でプローブにハイブリッド結合されるフルオレセイン標識標的から得られるフォトン数である。以下、完全対合プローブのハイブリッド形成強度を $I_{pm}$ で示し、不対合プローブのハイブリッド形成強度を $I_{mn}$ で示す。

【0114】ステップ954で、一対のプローブのハイブリッド形成強度を検索する。次のステップ956で、バックグラウンド信号強度を、その対の各ハイブリッド形成強度から差し引く。全ての生走査データから同時にバックグラウンドを差し引くようにしてもよい。

【0115】ステップ958で、プローブ対のハイブリ

ッド形成強度を差分閾値（D）および比の閾値（R）と比較する。即ち、プローブ対のハイブリッド形成強度の差（ $I_{pm} - I_{mn}$ ）が差分閾値以上であるかどうか、且つ、プローブ対のハイブリッド形成強度の比（ $I_{pm} / I_{mn}$ ）が比の閾値以上であるかどうかを判定される。これらの閾値は、通常、一つあるいは複数の遺伝子の発現状態を精度よく監視できるようにユーザーが定義した値である。例えば、差分閾値を20に比の閾値を1.2に設定するようにしてもよい。

【0116】 $I_{pm} - I_{mn} > D$ 、及び、 $I_{pm} / I_{mn} > R$ の場合には、ステップ960でNPOSを1増加させる。NPOSは、遺伝子が発現する可能性が高いことを示すハイブリッド形成強度を有するプローブ対の数を示す値である。即ち、遺伝子の発現状態の決定に、NPOSが用いられる。

【0117】ステップ962では、 $I_{mn} - I_{pm} > D$ 及び $I_{mn} / I_{pm} > R$ が成立するかどうかの判定がなされる。この式が成り立つ場合には、ステップ964でNNEGを1増加させる。NNEGは、遺伝子が発現しない可能性が高いことを示すハイブリッド形成強度を有するプローブ対の数を示す値である。NPOSと同様に、NNEGも遺伝子の発現状態の決定に用いられる。

【0118】遺伝子が発現していることを示す、あるいは、発現していないことを示すハイブリッド形成強度を有する各プローブ対に関して、ステップ966で、比の対数値（LR）と強度差値（IDIF）を計算する。LRは、プローブ対のハイブリッド形成強度の比（ $I_{pm} / I_{mn}$ ）の対数をとることにより求められる。また、IDIFは、プローブ対のハイブリッド形成強度の差（ $I_{pm} - I_{mn}$ ）をとることにより求められる。ステップ968で次のプローブ対のハイブリッド形成強度が存在すると判定されれば、処理をステップ954に戻して、検索を行う。

【0119】ステップ972で、決定行列を用いて、遺伝子が発現しているかどうかを判定する。決定行列では、N、NPOS、NNEG、LR（複数のLR）の値が用いられ、次の計算が行われる。

【0120】 $P1 = NPOS / NNEG$   
 $P2 = NPOS / N$   
 $P3 = (10 * \text{SUM}(LR)) / (NPOS + NNEG)$

【0121】これらのP値を用いて、遺伝子が発現しているかどうかを判定する。

【0122】説明のために、P値をいくつかの範囲に分割する。P1が2.1以上であれば、Aが真である。P1が2.1未満で1.8以上であれば、Bが真である。それ以外の場合には、Cが真である。即ち、P1は、A、B、Cの3つの範囲に分割される。これは、発明の理解を深めることを目的としてなされるものである。このようにして、全てのP値を以下のようにいくつかの範

10

20

30

40

50

図に分割する。

【0123】 $A = (P1 \geq 2, 1)$

$B = (2, 1 \geq P1 \geq 1, 8)$

$C = (P1 < 1, 8)$

【0124】 $X = (P2 \geq 0, 35)$

$Y = (0, 35 \geq P2 \geq 0, 20)$

$Z = (P2 < 0, 20)$

【0125】 $Q = (P3 \geq 1, 5)$

$R = (1, 5 \geq P3 \geq 1, 1)$

$S = (P3 < 1, 1)$

【0126】上記のブール値に従ってP値を所定の範囲に分割した後、遺伝子発現の判定を行う。

存在

$A \text{ and } (X \text{ or } Y) \text{ and } (Q \text{ or } R)$

$B \text{ and } X \text{ and } Q$

【0130】

どちらともいえない

$A \text{ and } X \text{ and } S$

$B \text{ and } X \text{ and } R$

$B \text{ and } Y \text{ and } (Q \text{ or } R)$

【0131】

不存在

それ以外の場合（例えば、Cの組み合わせ）

【0132】ステップ974で、ユーザーに結果が表示される。例えば、存在をP、どちらともいえないをM、不存在をAのように出力するようにしてもよい。

【0133】全てのブローブ対を処理して、遺伝子の発現状態を表示した後、ステップ975で、LRを10倍した値の平均値を計算する。さらに、NPOS及びNEGを増加させたブローブに関するIDIF値の平均を計算する。この平均値を、発現レベルを示す値として用いることもできる。また、これらの値を用いて、この実験と他の実験の定量的な比較を行うようにしてもよい。

【0134】ステップ976で定量的な測定を行う。一例として、現在の実験データを前の実験データと比較する（この場合、例えば、ステップ970で計算された値を用いる）。別の例として、生物検体中に存在する（バクテリア由来等の）既知量のRNAのハイブリッド形成強度と実験データを比較することもできる。このようにして、遺伝子発現状態の表示即ちコールが正しいかどうかを証明することもできるし、閾値を修正することもできるし、さらに、先に行われた実験の修正を行うことも可能である。

【0135】説明を簡単にするために、図20では一つの遺伝子に関する処理だけを説明したが、生物検体に含まれる複数の遺伝子の処理にこの方法を適用することもできる。従って、一つの遺伝子の解析に関する説明は、その方法を複数の遺伝子の処理に拡張できないことを示すものではない。

【0136】図21は、遺伝子発現監視ソフトウェアの画面表示（スクリーンディスプレイ）のレイアウト（構成）を示す。スクリーンディスプレイ1000は、グラフィックス表示領域1002とデータ表示領域1004

\*【0127】ここで、遺伝子の発現状態は、存在（発現）、どちらともいえない、不存在（発現しない）のいずれかで示される。式「 $A \text{ and } (X \text{ or } Y) \text{ and } (Q \text{ or } R)$ 」が成り立つ場合には、遺伝子が発現していると見なされる。言いかえると、 $P1 \geq 2, 1$ 、 $P2 \geq 0, 20$ 、 $P3 \geq 1, 1$ が成り立つ場合には、遺伝子が発現していると考えられる。また、式「 $B \text{ and } X \text{ and } Q$ 」が成り立つ場合にも、遺伝子が発現していると見なされる。

10 【0128】以上の説明に基づいて、遺伝子の発現状態を以下にまとめる。

\*【0129】

の2つの部分に分割される。グラフィックス表示領域には、ユーザーがデータを解釈する助けとなるグラフが表示される。また、データ表示領域には、ユーザーが遺伝子発現に関する基礎データを評価できるように、基礎データが表示される。

【0137】以下のスクリーンディスプレイの例に示すように、データ表示領域を縦欄と横欄で表示される表形式にすることが望ましい。各縦欄には、その欄に含まれるデータを示す見出しがつけられる。また、各横欄には、ある遺伝子に関する一つの実験あるいは実験の組み合わせから得られたデータが表示される。本明細書において、「実験」は、データを生成する工程を意味する。例えば、ハイブリッド結合されたチップの一つのイメージファイルから、多くの遺伝子に関する多くの「実験」データが生成される。また、実験は、異なったチップからデータを生成するものでもよい。

【0138】図22は、選択された遺伝子の解析結果を示すスクリーンディスプレイである。スクリーンディスプレイ1030のグラフィックス表示領域には、選択された遺伝子の各塩基位置における完全対合ブローブと不対合ブローブのハイブリッド形成強度が棒グラフの形で示されている。選択された遺伝子は、データ表示領域1034に反転表示されている。

【0139】データ表示領域には、複数の縦欄の見出しが表示されている。「Experiment Name」は、その実験に関してユーザーが定義した名称を意味する。「Gene Name」は、その遺伝子の名称を意味する。「Positive」及び「Negative」の数は、図20で説明したNPOS及びNEGの値を示す。「Pairs」は、その遺伝子の解析に用いられた完全符号ブローブと不対合ブローブのブ

ローブ対の数を示す。また、「Pos Fraction」は、遺伝子発現がポジティブであると判定されたプローブ対の割合（即ち、Positive/Pairs）を示す。

【0140】「Avg Ratio」欄は、ある遺伝子の全てのプローブに関する  $i_{pm}/i_{mn}$  の平均を示す。「Log A %」欄は、 $\log(i_{pm}/i_{mn})$  の平均を示す。「PM Excess」欄は、ユーザーが規定した閾値より大きなハイブリッド形成強度を示す完全対合プローブの数を示す。「NM Excess」欄は、ユーザーが規定した閾値より大きなハイブリッド形成強度を示す不對合プローブの数を示す。図23に示される「Pos/Neg」欄は、「Positive」欄と「Negative」欄の比を示す（「Negative」欄がゼロを含む場合には、「Inf」と表示される）。「Avg Diff」欄は、その遺伝子に関する平均強度差を示す。平均強度差は、図20のステップ975で計算された値である（即ち、(IDIF) 平均）。

【0141】「Abs Call」欄は、その実験に関する遺伝子発現コールを示す。この欄の値P、M、Aは、それぞれ、存在、どちらともいえない、不存在を表す。好適な実施例の遺伝子発現コールに関する詳細は、図20のステップ974を参照して既に説明した。

【0142】ユーザーが実験を選択すると、グラフィックス表示領域にユーザーがデータを解釈する助けとなるグラフが表示される。ユーザーは、ボタンバー1034を用いて、グラフィックス表示領域に表示されるべき一つあるいは複数のグラフを選択できる。さらに、所定の欄における値に従って、データ表示領域でデータを並べ変えることができる。

【0143】図24は、選択された遺伝子の解析結果を表す別のスクリーンディスプレイを示す。スクリーンディスプレイ1060のグラフィックス表示領域1062には、各塩基位置における完全対合プローブと不對合プローブのハイブリッド形成強度の比のグラフが表示されている。x軸は塩基位置を、y軸は、ハイブリッド形成強度の比を表す。また、グラフには統計学的な比の閾値がプロットされているが、この例では、閾値は1.2である。ユーザーは、このグラフを用いて、閾値より上及び下のプローブ対の数（ $i_{pm}/i_{mn}$ ）を解析することができる。このグラフには、さらに遺伝子名と実験名が表示されている。

【0144】図25は、選択された複数の遺伝子に関する実験結果の比較を表すスクリーンディスプレイを示す。スクリーンディスプレイ1160には、グラフィックス表示領域1062とデータ表示領域1164が含まれる。グラフィックス表示領域には、データ表示領域で選択された実験/遺伝子の各々に関して、各塩基位置における完全対合プローブと不對合プローブのハイブリッド形成強度の比のグラフが表示されている。好適な実施例では、選択された複数の遺伝子間の差異がはっきりとわかるように、実験名、遺伝子名及びデータプロットを

遺伝子ごとに異なった色で示す。

【0145】図26は、選択された複数の遺伝子に関する実験結果の比較を表す別のスクリーンディスプレイを示す。スクリーンディスプレイ1200のグラフィックス表示領域1202には、データ表示領域1204で選択された遺伝子の発現レベルが表示されている。選択された遺伝子の発現レベルは、棒グラフの形で示されている。好適な実施例では、発現レベルは、平均強度差として表される（図20の(IDIF) 平均を参照）。このグラフには、さらに遺伝子名と実験名が表示されている。

【0146】図27は、選択された複数の遺伝子に関する実験結果の比較を、グラフィックス表示領域に、いくつかのグラフの形で表示する、さらに別のスクリーンディスプレイを示す。スクリーンディスプレイ1230のグラフィックス表示領域1232には、データ表示領域1234で選択された遺伝子を解析した結果のグラフがいくつか示されている。この例では、選択された遺伝子に関する発現レベルグラフ1236、平均強度差グラフ1238、及び、ハイブリッド形成強度グラフ1240が表示されている。

【0147】図28及び図29は、基準走査データと実験走査データとを比較することにより、遺伝子の発現状態を判定する工程を示すフローチャートである。例えば、遺伝子が発現していることがわかっている生物検体から基準走査データをとるようにしてもよい。即ち、この走査データを別の生物検体と比較することにより、その遺伝子が発現しているかどうかを判定する。また、有機生命体において遺伝子の発現が時間とともにどのように変化するかを判定することもできる。ここで、「基準（基準）」は、標準点として用いられるものであることを意味する。

【0148】ステップ1302で、コンピュータシステムに、基準から得られたN対の完全対合プローブと不對合プローブの生の走査データが入力される。基準から得られた完全対合プローブのハイブリッド形成強度を  $i_{pm}$  で表し、同じく基準から得られた不對合プローブのハイブリッド形成強度を  $i_{mn}$  で表す。ステップ1304で、バックグラウンドの信号強度を各対の基準走査データのハイブリッド形成強度から差し引く。

【0149】一方、ステップ1306で、コンピュータシステムに、実験の標的である生物検体から得られたN対の完全対合プローブと不對合プローブの生の走査データが入力される。実験検体から得られた完全対合プローブのハイブリッド形成強度を  $j_{pm}$  で表し、同じく実験検体から得られた不對合プローブのハイブリッド形成強度を  $j_{mn}$  で表す。ステップ1308で、バックグラウンドの信号強度を各対の実験走査データのハイブリッド形成強度から差し引く。

【0150】次に、ステップ1310で、 $i-j$  対のハ



イブリッド形成強度を正規化する。例えば、 $i-j$ 対のハイブリッド形成強度を、対照ブロープのハイブリッド形成強度で割るようにしてもよい。

【0151】ステップ1312で、ブロープの $i-j$ 対のハイブリッド形成強度を差分閾値(DDIF)および比の閾値(RDIF)と比較する。具体的には、ある対( $J_{pm}-J_{mn}$ )のハイブリッド形成強度と別の対( $I_{pm}-I_{mn}$ )のハイブリッド形成強度の差が差分閾値以上であるかどうか。また、ある対( $J_{pm}-J_{mn}$ )のハイブリッド形成強度と別の対( $I_{pm}-I_{mn}$ )のハイブリッド形成強度の比が比の閾値以上であるかどうかの判定がなされる。これらの閾値は、通常、一つあるいは複数の遺伝子の発現監視を精度よく行えるようにユーザーが規定した値である。

【0152】( $J_{pm}-J_{mn})-(I_{pm}-I_{mn}) \geq DDIF$ 、且つ、( $J_{pm}-J_{mn})/(I_{pm}-I_{mn}) \geq RDIF$ が成立する場合には、ステップ1314でNINCを1増加させる。一般に、NINCは、実験検体のブロープ対の遺伝子発現が基準検体の遺伝子発現よりも大きい(即ち、増加している)可能性が高いことを示す値である。NINCを用いて、基準検体と比較して実験検体の遺伝子発現が大きい(即ち、増加している)か、小さい(即ち、減少している)か、あるいは、変化なしかの判定がなされる。

【0153】ステップ1316で、( $J_{pm}-J_{mn})-(I_{pm}-I_{mn}) \geq DDIF$ 、且つ、( $J_{pm}-J_{mn})/(I_{pm}-I_{mn}) \geq RDIF$ が成立するかどうかの判定がなされる。この式が成り立つ場合には、NDECを1増加させる。一般に、NDECは、実験検体のブロープ対の遺伝子発現が基準検体の遺伝子発現よりも小さい(即ち、減少している)可能性が高いことを示す値である。NDECを用いて、基準検体と比較して実験検体の遺伝子発現が大きい(即ち、増加している)か、小さい(即ち、減少している)か、あるいは、変化なしかの判定がなされる。

【0154】実験検体の遺伝子発現が大きいか、あるいは、小さいかを示すハイブリッド形成強度を有する各対のブロープに関して、NPOS、NNEG、及びLR値を計算する。これらの値の計算に関しては、図20を参照して前述した。それぞれの値の後ろにBあるいはEをつけて、その値が基準検体のものか実験検体のものかを示す。ステップ1322でハイブリッド形成強度を有する別のブロープ対が存在すると判定された場合には、上述と同様の方法で処理を行う。

【0155】次に、処理は図29に移って、ステップ1324で、基準検体と実験検体の両方に関して、絶対決\*

\*定演算を実行する。絶対決定演算は、基準検体と実験検体の各々に関して、遺伝子が発現しているか、どちらともいえないか、あるいは、発現していないか、を示すものである。従って、好適な実施例では、各検体に関して図20のステップ972と974を実行することにより、このステップの処理が実現される。この処理を行うことによって、各検体の遺伝子発現状態が示される。

【0156】次のステップ1326では、決定行列を用いて、2つの検体間の遺伝子発現の差を求める。この決定行列では、上記のように計算された値N、NPOS、NPOSE、NNEG、NNEGB、NNEG、NINC、NDEC、LRB、LREが用いられる。決定行列の計算は、NINCがNDEC以上であるかどうかによって変わる。即ち、以下のような計算が行われる。

【0157】 $NINC \geq NDEC$ の場合には、次の4つのP値が求められる。

【0158】 $P1 = NINC / NDEC$

$P2 = NINC / N$

$P3 = ((NPOSE - NPOS) - (NNEG - NNEGB)) / N$

$P4 = 10 * SUM(LRE - LRB) / N$

【0159】これらのP値を用いて、2つの検体間の遺伝子発現の差を求める。

【0160】上記と同様に、説明のために、P値を所定の範囲に分割する。P値は、全て、以下のような範囲に分割される。

【0161】 $A = (P1 \geq 2.8)$

$B = (2.8 > P1 \geq 2.0)$

$C = (P1 < 2.0)$

【0162】 $X = (P2 \geq 0.34)$

$Y = (0.34 > P2 \geq 0.24)$

$Z = (P2 < 0.24)$

【0163】 $M = (P3 \geq 0.20)$

$N = (0.20 > P3 \geq 0.12)$

$O = (P3 < 0.12)$

【0164】 $Q = (P4 \geq 0.9)$

$R = (0.9 > P4 \geq 0.5)$

$S = (P4 < 0.5)$

【0165】上記のブール値に従ってP値を所定の範囲に分割した後、2つの検体間の遺伝子発現の差を求める。

【0166】この場合、 $NINC \geq NDEC$ であるため、遺伝子発現の変化は、増加、増加傾向、あるいは、変化なしで表される。以下に、遺伝子発現状態をまとめて示す。

【0167】

増加

A and (X or Y) and (Q or R) and (M or N or O)

A and (X or Y) and (Q or R or S) and (M or N)

B and (X or Y) and (Q or R) and (M or N)

A and X and (Q or R or S) and (M or N or O)



【0168】

増加傾向 A or Y or S or O  
 B and (X or Y) and (Q or R) and O  
 B and (X or Y) and S and (M or N)  
 C and (X or Y) and (Q or R) and (M or N)

【0169】

変化なし それ以外の全ての場合（例えば、あらゆる2の組み合わせ）

【0170】ユーザーへの出力は、増加をI、増加傾向をMI、変化なしをNCのように示してもよい。

【0171】 $NINC < NDEC$ の場合には、次の4つ10のP値が求められる。

【0172】 $P1 = NDEC / NINC$

$P2 = NDEC / N$

$P3 = ((NNEG - NNEGB) - (NPOSE - NPOSB)) / N$

$P4 = 10 * SUM(LRE - LRB) / N$

【0173】これらのP値を用いて、2つの検体間の遺伝子発現の差を求める。

減少 A and (X or Y) and (Q or R) and (M or N or O)  
 A and (X or Y) and (Q or R or S) and (M or N)  
 B and (X or Y) and (Q or R) and (M or N)  
 A and X and (Q or R or S) and (M or N or O)

【0177】

減少傾向 A or Y or S or O  
 B and (X or Y) and (Q or R) and O  
 B and (X or Y) and S and (M or N)  
 C and (X or Y) and (Q or R) and (M or N)

【0178】

変化なし それ以外の全ての場合（例えば、あらゆる2の組み合わせ）

【0179】ユーザーへの出力は、減少をD、減少傾向をDI、変化なしをNCのように示してもよい。

【0180】このようにして、基準検体と実験検体との間の遺伝子発現の相対的な差を求めることができる。

（例えば、ステップ1324で）両方の検体において遺伝子が発現状態であると表示され、且つ、以下の式が全て成立する場合には、追加試験を行って、I、MI、D、M、D（即ちNC以外の）コールをNCに変更するように構成してもよい。

【0181】 $(IDIFB)の平均 \geq 200$

$(IDIFE)の平均 \geq 200$

1.  $4 = \{(IDIFE)の平均\} / \{(IDIFB)の平均\} \geq 0.7$

【0182】即ち、両方の検体で遺伝子が発現した場合に、各検体の平均強度差が比較的大きいか、あるいは、両方の検体の平均強度差がほとんど同じであれば、増加あるいは減少のコールは（それが傾向に留まっているかどうかに関わらず）変化なしコールに変更される。IDIFB及びIDIFEは、各検体に関するIDIFの総和をNで割ることにより計算される。

【0183】次に、ステップ1328で、定量的な差の

\* 【0174】これらのP値は、上記の $NINC \geq NDEC$ の場合と同じいくつかの範囲に分割する。P値は上述と同じ範囲に分割されるので、ここでは簡略化のために説明を繰り返さない。但し、以下に説明するように、これらの範囲は、通常、2つの検体間の遺伝子発現の異なった変化状態を示す。

【0175】この場合、 $NINC < NDEC$ であるため、遺伝子発現の変化は、減少、減少傾向、あるいは、変化なしで表される。以下に、遺伝子発現状態をまとめて示す。

\* 【0176】

30 評価を行うための値を計算する。各対に関して、 $((J_{pm} - J_{mm}) - (I_{pm} - I_{mm}))$ の平均を計算する。また、 $J_{pm} - J_{mm}$ の平均と $I_{pm} - I_{mm}$ の平均の比を計算する。ステップ1330で、これらの値を用いて、他の実験結果との比較を行うことができる。

【0184】図30は、2つの実験における遺伝子発現の変化を表すスクリーンディスプレイを示す。スクリーンディスプレイ1400には、グラフィックス表示領域1402とデータ表示領域1404が含まれる。ユーザーがある所定の遺伝子に関して2つの実験を選択することにより、その実験結果の比較が行われる。ここでは、説明をわかりやすくするために、基準データと実験データの2つを呼び出して、それらのデータの比較を行う場合を考える。例えば、ユーザーは、「8182506」と名付けられた遺伝子に関して2つの実験を選択する。2つの実験結果の比較自体が実験であるため、ユーザーは、図30のデータ表示領域に「foo」と表示されるように、実験名を入力することができる。図31は、2つの実験における遺伝子発現の変化を表す別のスクリーンディスプレイを示す。

【0185】システムは、図32で説明した方法に従っ

て、選択された2つの実験における遺伝子発現の変化を求める。データ表示領域には、この比較によって生成されたデータが表形式で示されている。「Experiment Name」は、比較実験に対してユーザーが規定した名称を示す。「Gene Name」は、遺伝子の名称を示す。「Inc」及び「Dec」の数値は、図28を参照して説明したように、NINCおよびNDICの値を示す。具体的には、「Inc」は、完全対合プローブと不対合プローブのハイブリッド形成強度の差及び比が実験データにおいて有意に大きいような遺伝子における塩基位置の数を示す。

【0186】「Inc Ratio」欄は、ハイブリッド形成強度が増加した塩基位置の数を解析対象の遺伝子における全ての塩基位置の数で割った値を示す。「Dec Ratio」欄は、ハイブリッド形成強度が減少した塩基位置の数を解析対象の遺伝子における全ての塩基位置の数で割った値を示す。「Pos Change」欄は、実験データと基準データにおいて正の得点を有するプローブ対の数の差を示す。「Neg Change」欄は、実験データと基準データにおいて負の得点を有するプローブ対（完全対合プローブと不対合プローブ）の数の差を示す。

【0187】「Inc/Dec」欄は、実験データにおいてハイブリッド形成強度が増加したプローブ対の数とハイブリッド形成強度が減少したプローブ対の数の比を示す。「Avg Diff」欄は、実験データにおける平均強度差を示す。

【0188】（図示されていない）「Diff Call」欄は、所定の遺伝子に関する2つの実験における発現レベルの変化を示す。この欄において、Iは遺伝子発現の増加を、Mは遺伝子発現の増加傾向を、Dは遺伝子発現の減少を、MDは遺伝子発現の減少傾向を、NCは変化なしを、?は未知であることを、それぞれ示す。好適な実施例において、発現レベルの変化は、図29のステップ1326を参照して説明したように計算される。

【0189】遺伝子発現の変化の演算に加えて、ユーザーは、データ解析のためにグラフを選択することもできる。グラフィックス表示領域1402には、基準データと実験データを示す3つのグラフが表示されている。

【0190】図32は、2つの実験における遺伝子発現の変化を3次元棒グラフの形で示したスクリーンディスプレイを示す。スクリーンディスプレイ1440のグラフィックス表示領域1442には、データ表示領域1444で選択された遺伝子の発現レベルを示す3次元棒グラフが表示されている。ユーザーは、一つあるいは複数の遺伝子をデータ表示領域で選択して、これら遺伝子の発現レベルを示す3次元棒グラフを作成するようにシステムに命令する。好適な実施例において、発現レベルは、平均強度差（即ち、（IDIF）の平均）として与えられる。3次元棒グラフにより、ユーザーは、簡単に、複数の遺伝子の発現レベルを調べることができる。また、複数の実験結果から選択された、同じ遺伝子に関

するデータを同時に示して回転させることにより、発現レベルの差を表示するようにすることもできる。

#### 【0191】結語

上述の説明は例示を目的としたもので、何ら発明を限定するものではない。当業者に明らかなように、この開示に基づき、本発明を様々な変更することができる。例えば、本発明では、（天然のものも人工のものも含めて）DNAの評価に関して説明しているが、本発明の方法を、RNAのような他の物質が合成されたチップの解析に適用することも可能である。従って、本発明の範囲は、先の説明によって決まるものではなく、特許請求の範囲およびそれと等価と考えられる全範囲によって決定される。

#### 【図面の簡単な説明】

【図1】本発明のソフトウェアを実行するために用いられるコンピュータシステムの例を示す図。

【図2】典型的なコンピュータシステムを示すシステムブロック図。

【図3】DNAやRNA等の生体物質の配列（アレイ）を形成し、解析するための全体的なシステムを示す図。

【図4】全体的なシステムに関するソフトウェアの実施例を示す図。

【図5】全体的なシステムで形成されたチップの全体的な構成を示す図。

【図6】チップ上の核酸プローブが標識された標的に結合する様子を概念的に示す図。

【図7】チップ上に列で配置された核酸プローブを示す図。

【図8】図7のような参照配列と共に、チップ上の標的のハイブリッド形成パターンを示す図。

【図9】標準構造及び代替構造を示す図。

【図10】チップから得られたハイブリッド形成強度を表すスクリーンディスプレイを示す図。

【図11】関連するプローブのハイブリッド形成強度から塩基コールを演算する方法を示すフローチャート。

【図12】関連するプローブのハイブリッド形成強度から塩基コールを演算する別の方法を示すフローチャート。

【図13】一つのユニットグループにおいて塩基コール処理を行う方法を示すフローチャート。

【図14】複数のユニットグループにおいて塩基コール処理を行う方法を示すフローチャート。

【図15】一つのユニットグループに関して一つの塩基を呼び出す方法を示すフローチャート。

【図16】塩基コール処理を実行するために最適なユニットグループを選択する方法を示すフローチャート。

【図17】一つあるいは複数のチップから得られた実験データに基づきスクリーンディスプレイを解析するためのスクリーンディスプレイを示す図。

【図18】一つあるいは複数のチップから得られた実験データに基づきヌクレオチドを解析するためのスクリーンディスプレイを示す図。

【図19】完全対合プローブと不対合プローブ対のハイブリッド形成強度を比較することにより、ある遺伝子の発現状態を監視する方法を示すフローチャート。

【図20】決定行列を用いて遺伝子が発現しているかどうかを判定する方法を示すフローチャート。

【図21】遺伝子発現を監視するソフトウェアのスクリーンディスプレイの配置を示す図。

【図22】選択された遺伝子の解析結果を表すスクリーンディスプレイを示す図。

【図23】選択された遺伝子の解析結果を表すスクリーンディスプレイを示す図。

【図24】選択された遺伝子の解析結果を表す別のスクリーンディスプレイを示す図。

【図25】選択された複数の遺伝子に関する実験結果の比較を表すスクリーンディスプレイを示す図。

\* 【図26】選択された複数の遺伝子に関する実験結果の比較を表す別のスクリーンディスプレイを示す図。

【図27】選択された複数の遺伝子に関する実験結果の比較をグラフィックス表示領域における複数のグラフを用いて表す別のスクリーンディスプレイを示す図。

【図28】基準走査データと実験走査データとを比較することにより、ある遺伝子の発現状態を決定する方法を示すフローチャート。

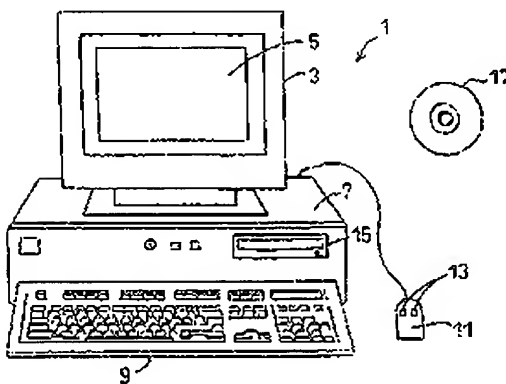
【図29】基準走査データと実験走査データとを比較することにより、ある遺伝子の発現状態を決定する方法を示すフローチャート。

【図30】2つの実験における遺伝子発現の変化を表すスクリーンディスプレイを示す図。

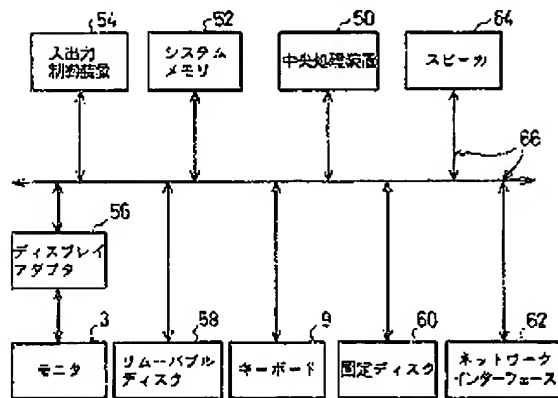
【図31】2つの実験における遺伝子発現の変化を表すスクリーンディスプレイを示す図。

【図32】2つの実験における遺伝子発現の変化を3次元棒グラフの形で表すスクリーンディスプレイを示す図。

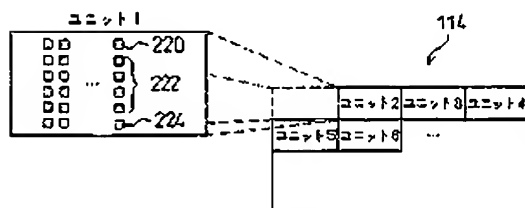
【図1】



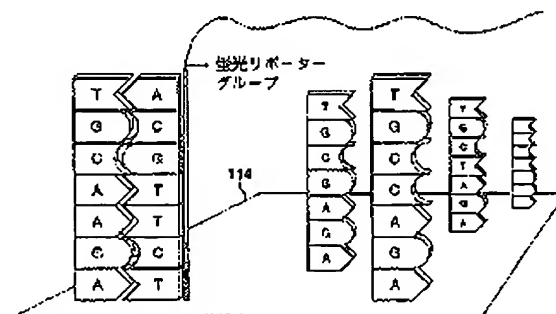
【図2】



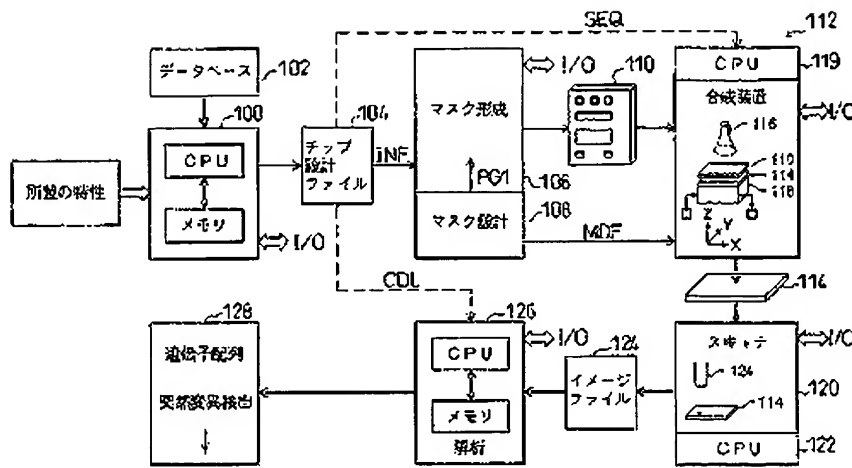
【図5】



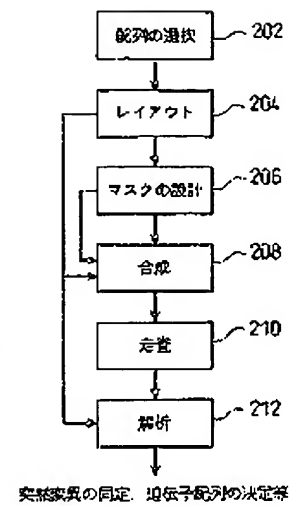
【図6】



【图3】

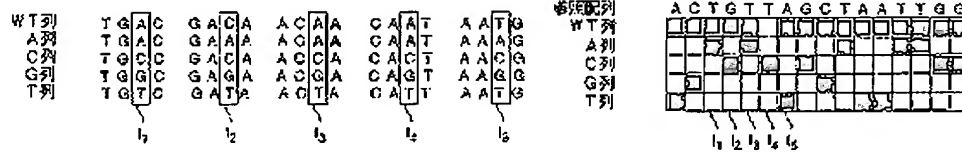


【圖4】

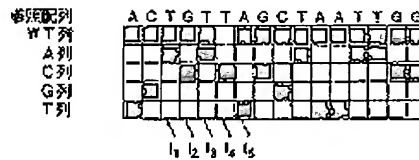


【圖 7】

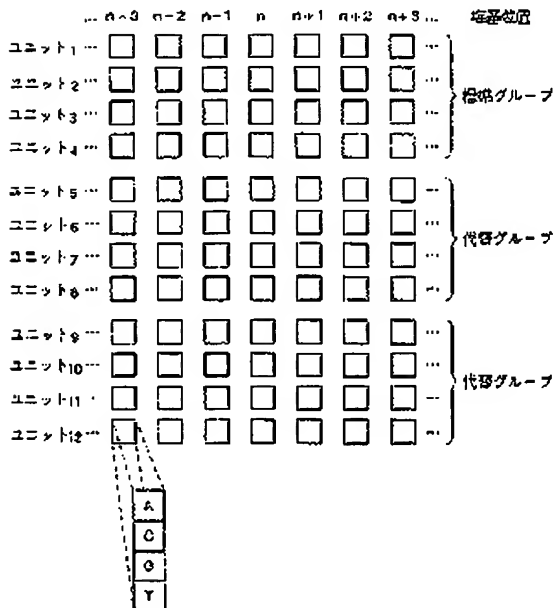
序列排列 A C T<sub>1</sub> G<sub>2</sub> T<sub>3</sub> T<sub>4</sub> A<sub>5</sub> G C T A A T T G G · 5'



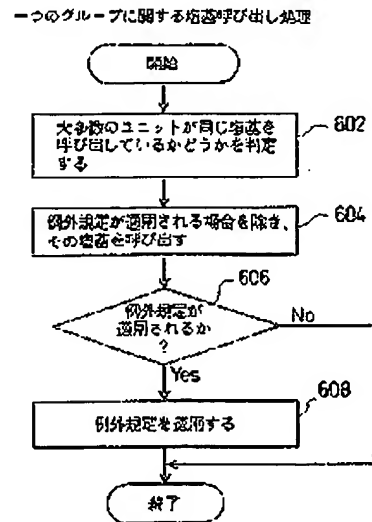
【图8】



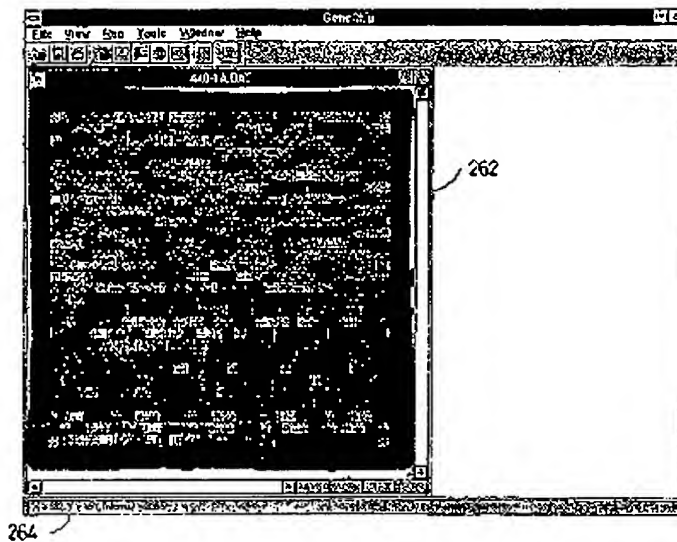
【図9】



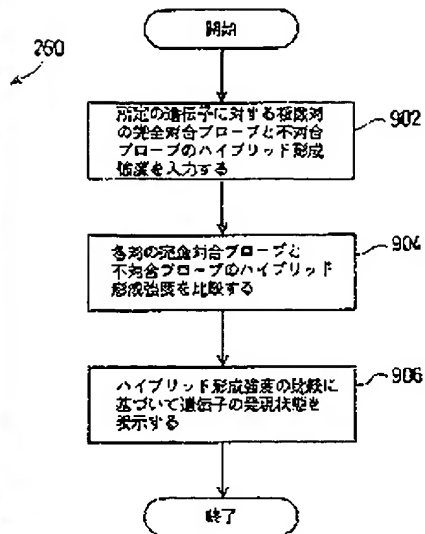
【图 15】



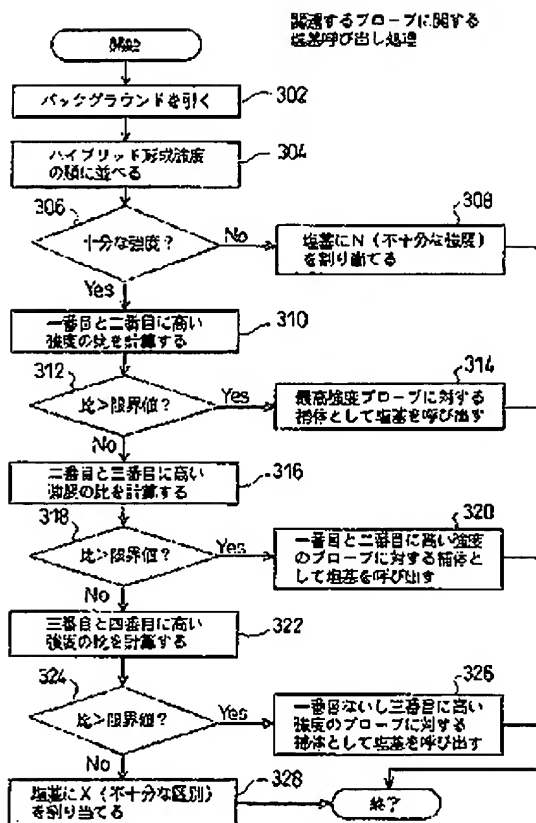
【図10】



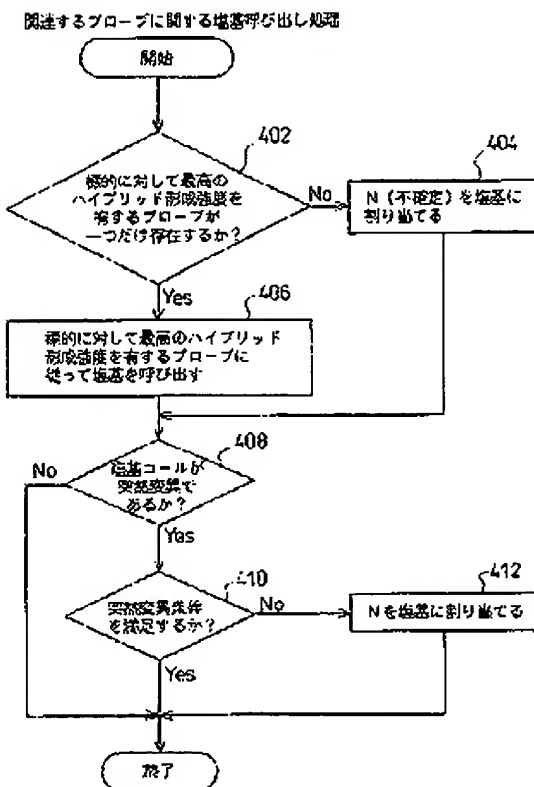
【図19】



【図11】

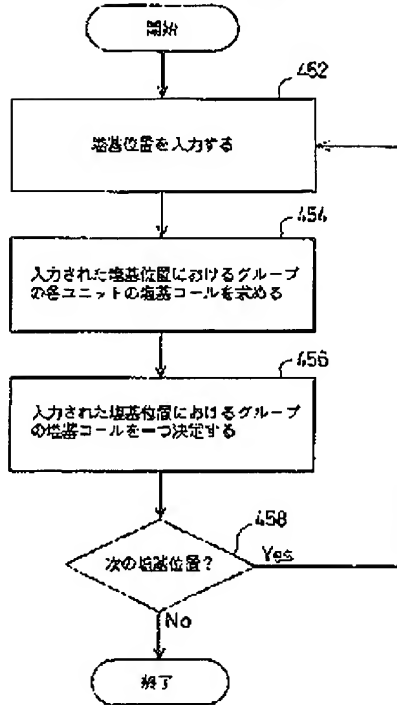


【図12】



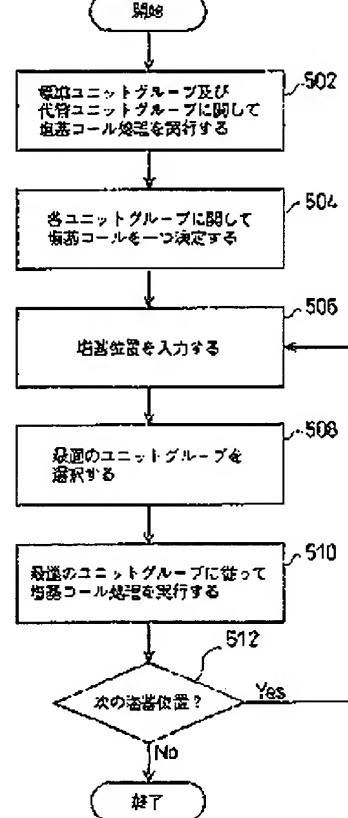
【図13】

一つのグループにおける塩基呼び出し処理

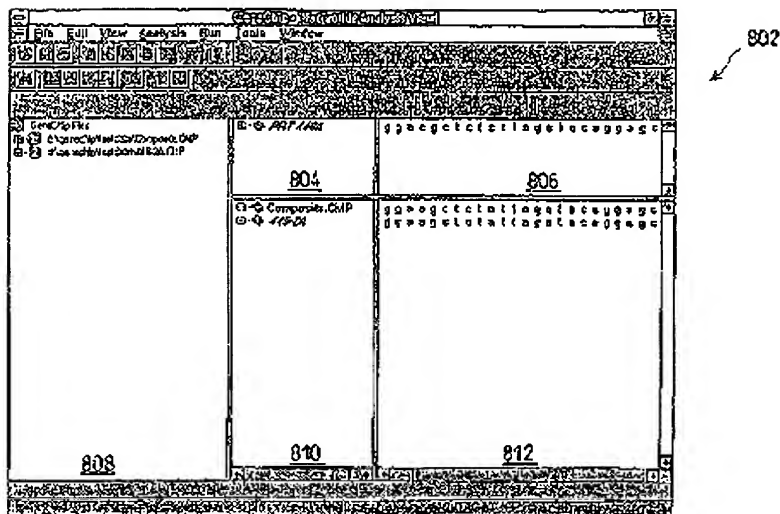


【図14】

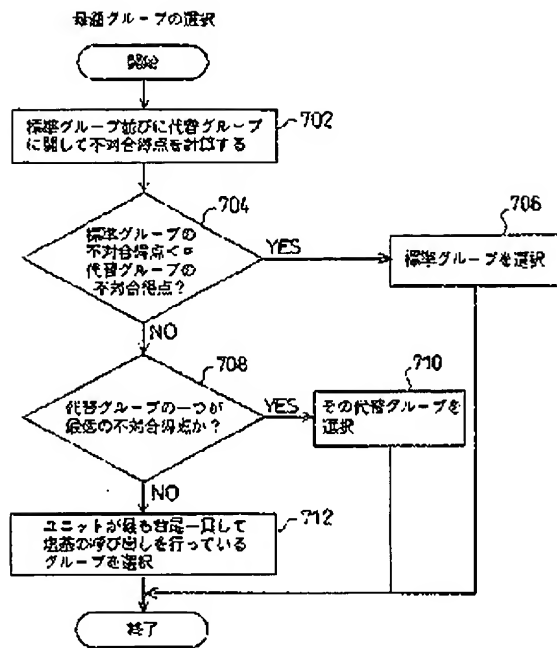
複数のグループに関する塩基呼び出し処理



【図17】



【圖 16】



```

graph TD
    A((A)) --> B[1324 基底及び副基底について  
総対決定演算を行う]
    B --> C[1326 決定行列を用いてコーネルを求める]
    C --> D[1328 
$$\frac{(J_{pn} - J_{qn})}{(J_{pn} - J_{qn}) \text{ 平均}} = \frac{(I_{pn} - I_{qn})}{(I_{pn} - I_{qn}) \text{ 平均}}$$
]
    D --> E[1330 定数か? 差の検出]
    E --> F([終了])
  
```

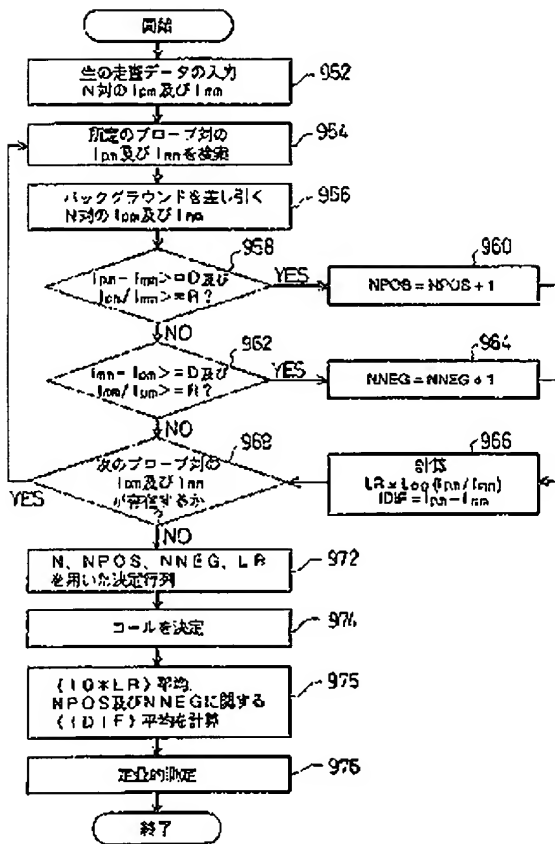
Flowchart illustrating the determination method for the number of rows in the matrix:

- Step 1324: Perform total determination calculation for the basis and sub-basis.
- Step 1326: Obtain the Cornell using the determination matrix.
- Step 1328: Calculate the ratio: 
$$\frac{(J_{pn} - J_{qn})}{(J_{pn} - J_{qn}) \text{ 平均}} = \frac{(I_{pn} - I_{qn})}{(I_{pn} - I_{qn}) \text{ 平均}}$$
- Step 1330: Check if the value is constant or if there is a difference.
- End (終了).

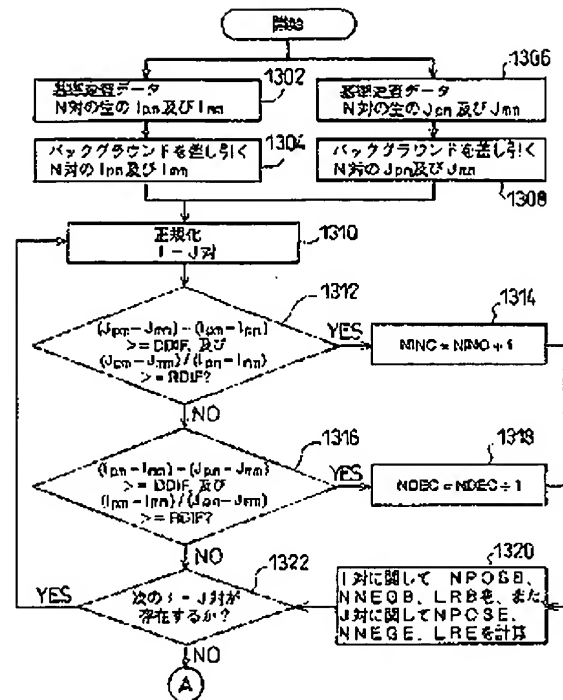
[illegible]

602

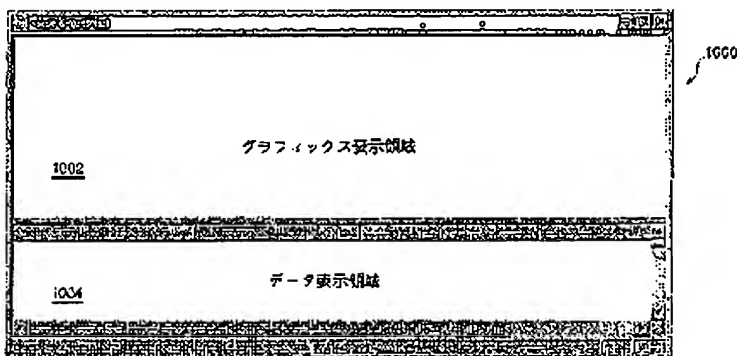
【図20】



【図28】



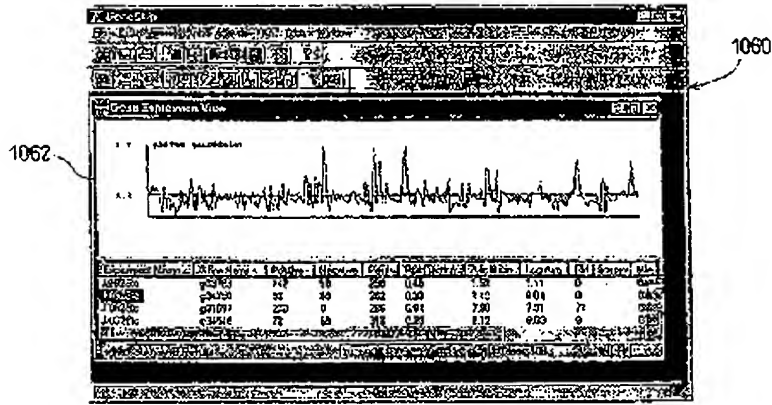
【図21】



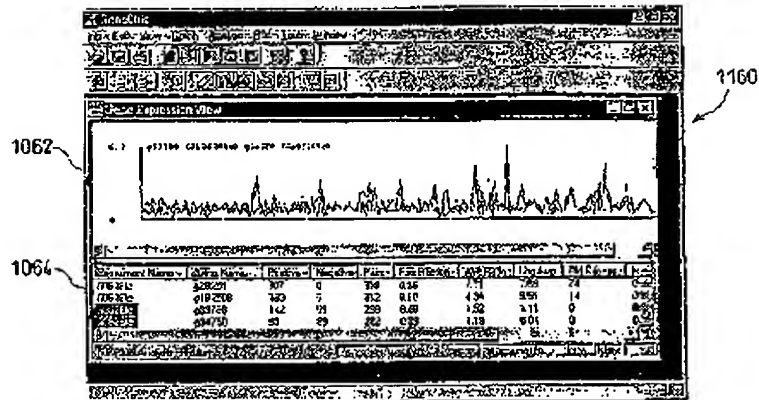


[illegible]

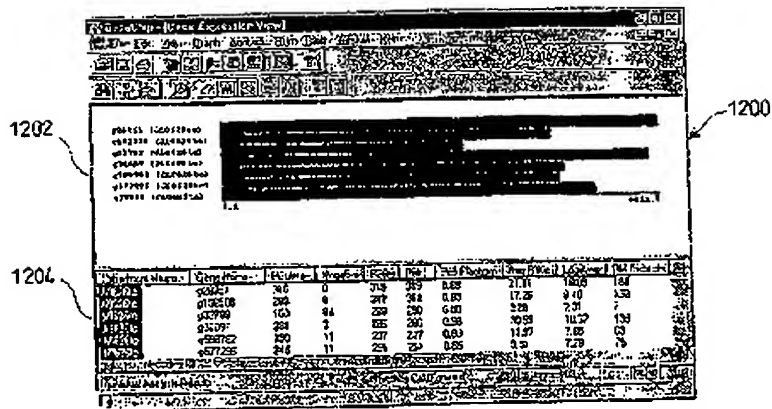
【图 24】



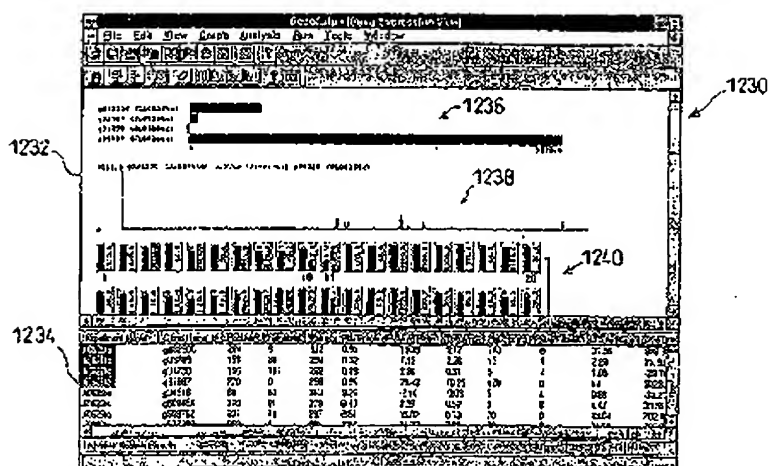
【图 25】



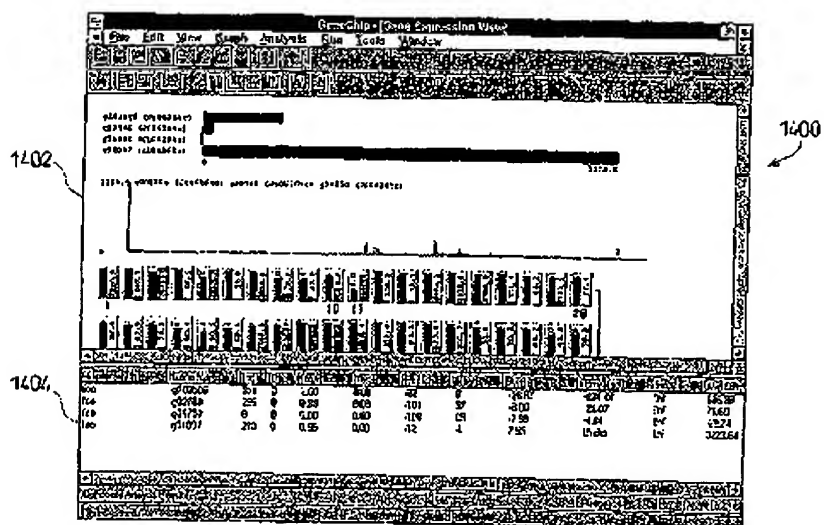
【圖26】



【図27】



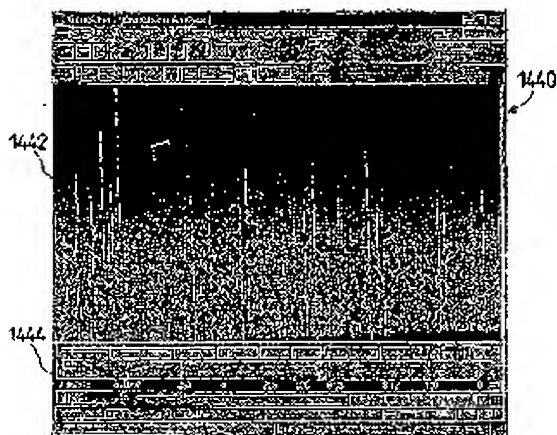
【図30】



【図31】

1404

【図32】



フロントページの続き

(72)発明者 モリス・エス・マクドナルド  
アメリカ合衆国 カリフォルニア州94305  
スタンフォード、ソノマ・テラス、843  
(72)発明者 マイケル・ビー・ミットマン  
アメリカ合衆国 カリフォルニア州94303  
バロ・アルト、セイント・フランシス・  
ドライブ、2377

(72)発明者 デビット・ジェイ・ロックハート  
アメリカ合衆国 カリフォルニア州94041  
マウンテン・ビュー、マウンテン・ビュ  
ー・アベニュー、610  
(72)発明者 ミン・スー・ホー  
アメリカ合衆国 カリフォルニア州95133  
サン・ホセ、フォックスリッジ・ウェ  
イ、922

(72)発明者 デレク・バーンハート  
アメリカ合衆国 カリフォルニア州94087  
サニーヴェール, シラス・ウェイ、422

(72)発明者 ルイス・シー・・ジェボンス  
アメリカ合衆国 カリフォルニア州94087  
サニーヴェール, ロマーナ・アベニュー、 701

## 〔外国語明細書〕

## 1 Title of Invention

## COMPUTER-AIDED TECHNIQUES FOR ANALYZING BIOLOGICAL SEQUENCES

## 2 Claims

1. In a computer system, a method of analyzing a sample nucleic acid sequence, the method comprising the steps of:

inputting a plurality of base calls for each base position along at least a portion of the sample nucleic acid sequence;

for each base position, analyzing the plurality of base calls to generate a single base call; and

displaying single base calls for base positions along the at least a portion of said sample nucleic acid sequence, each of the single base calls being derived from the plurality of base calls for a specific base position in the sample nucleic acid sequence.

2. The method of claim 1, wherein the analyzing step comprises the steps of:

for each base position, determining a base call of the plurality of base calls which occurs most often; and

generating the single base call as the base call that occurs most often at the base position.

3. The method of claim 1, further comprising the step of displaying a screen icon which when activated by a user causes the plurality of base calls at each base position to be displayed.

4. The method of claim 1, further comprising the step of displaying a screen icon which when activated by a user causes the plurality of base calls at each base position not to be displayed.

5. The method of claim 1, further comprising the step of displaying the plurality of base calls at each base aligned with the single base calls according to base position.

6. The method of claim 5, further comprising the step of displaying with each base call of the plurality of base calls hybridization intensities indicating hybridization affinity of a probe and the sample nucleic acid sequence,

整理番号 = P A 5 2 C 2 9 5

ページ (2)

wherein each base call is determined by an analysis of the hybridization intensities.

7. In a computer system, a method of calling an unknown base in a sample nucleic acid sequence, the method comprising the steps of:

receiving hybridization intensities for a plurality of sets of nucleic acid probes, each hybridization intensity indicating a hybridization affinity between a nucleic acid probe and the sample nucleic acid sequence;

computing a base call for the unknown base for each set of probes; and

computing a single base call for the plurality of sets of probes according to the base call for the unknown base which occurs most often for the plurality of sets of probes.

8. The method of claim 7, wherein each set of probes was generated according to a same reference sequence.

9. The method of claim 7, further comprising the step of checking exception rules that specify the single base call for the plurality of sets of nucleic acid probes under certain conditions.

10. In a computer system, a method of dynamically changing parameters for a computer-implemented base calling procedure, the method comprising the steps of:

generating base calls for at least a portion of a sample nucleic acid sequence utilizing the base calling procedure, the base calling procedure including a parameter that is changeable by a user;

displaying the base calls for the at least a portion of a sample nucleic acid sequence;

displaying the parameter of the base calling procedure;

receiving input from the user specifying a new value for the parameter of the base calling procedure;

generating updated base calls for the at least a portion of a sample nucleic acid sequence utilizing the base calling procedure and the new value for the parameter; and

displaying the updated base calls for the at least a portion of a sample nucleic acid sequence.

11. The method of claim 10, further comprising the step of displaying a plurality of user-changeable parameters for the base calling procedure.

整理番号 = F A 5 2 C 2 9 5

ページ (3)

12. The method of claim 10, wherein the parameter is selected from the group consisting of a constant, threshold, and range.

13. In a computer system, a method of monitoring expression of a gene in a sample nucleic acid sequence, the method comprising the steps of:

inputting a plurality of hybridization intensities of pairs of perfect match and mismatch probes, the perfect match probes being perfectly complementary to the gene and the mismatch probes having at least one base mismatch with the gene, and the hybridization intensities indicating hybridization intensity between the perfect match and mismatch probes and the sample nucleic acid sequence;

comparing the hybridization intensities of each pair of perfect match probes in order to generate a gene expression call of the sample nucleic acid sequence; and

displaying the gene expression call.

14. The method of claim 13, further comprising the step of comparing a difference between hybridization intensities of perfect match and mismatch probes at a base position to a difference threshold.

15. The method of claim 13, further comprising the step of comparing a quotient of hybridization intensities of perfect match and mismatch probes at a base position to a ratio threshold.

16. The method of claim 13, further comprising the step of utilizing a decision matrix to determine the gene expression call.

17. The method of claim 13, wherein the gene expression call is selected from the group consisting of expressed, marginal, and absent.

18. In a computer system, a method of monitoring expression of a gene in a sample nucleic acid sequence, the method comprising the steps of:

inputting a plurality of hybridization intensities of pairs of perfect match and mismatch probes, the perfect match probes being perfectly complementary to the gene and the mismatch probes having at least one base mismatch with the gene, and the hybridization intensities indicating hybridization intensity between the perfect match and mismatch probes and the sample nucleic acid sequence;

comparing the hybridization intensities of each pair of perfect match



整理番号 = P A 5 2 C 2 9 5

ページ (4)

probes; and

generating a gene expression call of the sample nucleic acid sequence.

19. The method of claim 18, further comprising the step of comparing a difference between hybridization intensities of perfect match and mismatch probes at a base position to a difference threshold.

20. The method of claim 18, further comprising the step of comparing a quotient of hybridization intensities of perfect match and mismatch probes at a base position to a ratio threshold.

21. The method of claim 18, further comprising the step of utilizing a decision matrix to determine the gene expression call.

22. The method of claim 18, wherein the gene expression call is selected from the group consisting of expressed, marginal, and absent.

23. In a computer system, a method of monitoring change in expression of a gene in a sample nucleic acid sequence, the method comprising the steps of:

inputting a plurality of hybridization intensities of pairs of perfect match and mismatch probes, the perfect match probes being perfectly complementary to the gene and the mismatch probes having at least one base mismatch with the gene, and the hybridization intensities indicating hybridization intensity between the perfect match and mismatch probes and the sample nucleic acid sequence;

comparing the hybridization intensities of each pair of perfect match probes in order to generate a gene expression level of the sample nucleic acid sequence;

determining a change in expression by comparing the gene expression level to a baseline gene expression level; and

displaying the change in expression of the gene in the sample nucleic acid.

24. The method of claim 23, wherein the change in expression is displayed as a graph.

25. The method of claim 23, further comprising the step of generating the baseline expression level according to the inputting and comparing

整理番号 = P A 5 2 C 2 9 5

ページ (5)

steps of claim 23.

26. The method of claim 23, further comprising the step of comparing hybridization intensities of perfect match and mismatch probes hybridizing with the sample nucleic acid sequence and hybridization intensities of perfect match and mismatch probes hybridizing with a baseline sequence to a difference threshold.

27. The method of claim 23, further comprising the step of comparing hybridization intensities of perfect match and mismatch probes hybridizing with the sample nucleic acid sequence and hybridization intensities of perfect match and mismatch probes hybridizing with a baseline sequence to a ratio threshold.

28. The method of claim 23, further comprising the step of utilizing a decision matrix to determine the change in expression of the gene in the sample nucleic acid.

29. The method of claim 23, wherein the change in expression of the gene in the sample nucleic acid is selected from the group consisting of increased, marginal increase, decreased, marginal decrease, and no change.

30. In a computer system, a method of monitoring change in expression of a gene in a sample nucleic acid sequence, the method comprising the steps of:

inputting a plurality of hybridization intensities of pairs of perfect match and mismatch probes, the perfect match probes being perfectly complementary to the gene and the mismatch probes having at least one base mismatch with the gene, and the hybridization intensities indicating hybridization intensity between the perfect match and mismatch probes and the sample nucleic acid sequence;

comparing the hybridization intensities of each pair of perfect match probes in order to generate a gene expression level of the sample nucleic acid sequence; and

determining a change in expression by comparing the gene expression level to a baseline gene expression level.

31. The method of claim 30, further comprising the step of generating the baseline expression level according to the inputting and comparing

整理番号=P A 5 2 C 2 9 5

ページ (8)

steps of claim 30.

32. The method of claim 30, further comprising the step of comparing hybridization intensities of perfect match and mismatch probes hybridizing with the sample nucleic acid sequence and hybridization intensities of perfect match and mismatch probes hybridizing with a baseline sequence to a difference threshold.

33. The method of claim 30, further comprising the step of comparing hybridization intensities of perfect match and mismatch probes hybridizing with the sample nucleic acid sequence and hybridization intensities of perfect match and mismatch probes hybridizing with a baseline sequence to a ratio threshold.

34. The method of claim 30, further comprising the step of utilizing a decision matrix to determine the change in expression of the gene in the sample nucleic acid.

35. The method of claim 30, wherein the change in expression of the gene in the sample nucleic acid is selected from the group consisting of increased, marginal increase, decreased, marginal decrease, and no change.

### 3 Detailed Description of Invention

The present invention relates to the field of computer systems. More specifically, the present invention relates to computer systems for analyzing biological sequences such as nucleic acid sequences.

Devices and computer systems for forming and using arrays of materials on a substrate are known. For example, PCT application WO92/10588, incorporated herein by reference for all purposes, describes techniques for sequencing or sequence checking nucleic acids and other materials. Arrays for performing these operations may be formed in arrays according to the methods of, for example, the pioneering techniques disclosed in U.S. Patent No. 5,143,854 and U.S. Patent Application No. 08/249,188, both incorporated herein by reference for all purposes.

According to one aspect of the techniques described therein, an array of nucleic acid probes is fabricated at known locations on a substrate or chip. A

整理番号 = P A 5 2 〇 2 9 5

ページ (7)

fluorescently labeled nucleic acid is then brought into contact with the chip and a scanner generates an image file (which is processed into a cell file) indicating the locations where the labeled nucleic acids bound to the chip. Based upon the cell file and identities of the probes at specific locations, it becomes possible to extract information such as the monomer sequence of DNA or RNA. Such systems have been used to form, for example, arrays of DNA that may be used to study and detect mutations relevant to cystic fibrosis, the P53 gene (relevant to certain cancers), HIV, and other genetic characteristics.

Innovative computer-aided techniques for base calling are disclosed in U.S. Patent Application Nos. 08/531,137 (attorney docket no. 16528X-008210), 08/528,655 (attorney docket no. 16528X-017600), and 08/618,834 (attorney docket no. 16528X-016400), which are all hereby incorporated by reference for all purposes. However, improved computer systems and methods are still needed to evaluate, analyze, and process the vast amount of information now used and made available by these pioneering technologies.

Additionally, there is a need for improved computer-aided techniques for monitoring gene expression. Many disease states are characterized by differences in the expression levels of various genes either through changes in the copy number of the genetic DNA or through changes in levels of transcription (e.g., through control of initiation, provision of RNA precursors, RNA processing, etc.) of particular genes. For example, losses and gains of genetic material play an important role in malignant transformation and progression. Furthermore, changes in the expression (transcription) levels of particular genes (e.g., oncogenes or tumor suppressors), serve as signposts for the presence and progression of various cancers.

Similarly, control of the cell cycle and cell development, as well as diseases, are characterized by the variations in the transcription levels of particular genes. Thus, for example, a viral infection is often characterized by the elevated expression of genes of the particular virus. For example, outbreaks of *Herpes simplex*, Epstein-Barr virus infections (e.g., infectious mononucleosis), cytomegalovirus, Varicella-zoster virus infections, parvovirus infections, human papillomavirus infections, etc. are all characterized by elevated expression of

整理番号=P A 5 2 C 2 9 5

ページ (8)

various genes present in the respective virus. Detection of elevated expression levels of characteristic viral genes provides an effective diagnostic of the disease state. In particular, viruses such as herpes simplex, enter quiescent states for periods of time only to erupt in brief periods of rapid replication. Detection of expression levels of characteristic viral genes allows detection of such active proliferative (and presumably infective) states.

#### SUMMARY OF THE INVENTION

The present invention provides innovative systems and methods for analyzing biological sequences such as nucleic acid sequences. The computer system may analyze hybridization intensities indicating hybridization affinity between nucleic acid probes and a sample nucleic acid sequence in order to call bases in the sample sequence. Multiple base calls may be combined to form a single base call. Additionally, the computer system may analyze hybridization intensities in order to monitor gene expression or the change in gene expression as compared to a baseline.

According to one aspect of the invention, a computer-implemented method of calling an unknown base in a sample nucleic acid sequence comprises the steps of: receiving hybridization intensities for a plurality of sets of nucleic acid probes, each hybridization intensity indicating a hybridization affinity between a nucleic acid probe and the sample nucleic acid sequence; computing a base call for the unknown base for each set of probes; and computing a single base call for the plurality of sets of probes according to the base call for the unknown base which occurs most often for the plurality of sets of probes. Typically, the single base call is displayed on a screen display and a user is afforded the opportunity to display or not display the base cases from which the single base call is derived.

According to another aspect of the invention, a method of dynamically changing parameters for a computer-implemented base calling procedure comprises the steps of: generating base calls for at least a portion of a sample nucleic acid sequence utilizing the base calling procedure, the base calling procedure including a parameter that is changeable by a user; displaying the base calls for the at least a portion of a sample nucleic acid sequence; displaying the parameter of the base

整理番号=PA520295

ページ (9)

calling procedure; receiving input from the user specifying a new value for the parameter of the base calling procedure; generating updated base calls for the at least a portion of a sample nucleic acid sequence utilizing the base calling procedure and the new value for the parameter; and displaying the updated base calls for the at least a portion of a sample nucleic acid sequence. Typically the user-changeable parameter is a constant, threshold, or range.

According to another aspect of the invention, a computer-implemented method of monitoring expression of a gene in a sample nucleic acid sequence comprises the steps of: inputting a plurality of hybridization intensities of pairs of perfect match and mismatch probes, the perfect match probes being perfectly complementary to the gene and the mismatch probes having at least one base mismatch with the gene, and the hybridization intensities indicating hybridization intensity between the perfect match and mismatch probes and the sample nucleic acid sequence; comparing the hybridization intensities of each pair of perfect match probes; and generating a gene expression call of the sample nucleic acid sequence. In preferred embodiments, the expression call is denoted as expressed, marginal, or absent.

According to another aspect of the invention, a computer-implemented method of monitoring change in expression of a gene in a sample nucleic acid sequence comprises the steps of: inputting a plurality of hybridization intensities of pairs of perfect match and mismatch probes, the perfect match probes being perfectly complementary to the gene and the mismatch probes having at least one base mismatch with the gene, and the hybridization intensities indicating hybridization intensity between the perfect match and mismatch probes and the sample nucleic acid sequence; comparing the hybridization intensities of each pair of perfect match probes in order to generate a gene expression level of the sample nucleic acid sequence; and determining a change in expression by comparing the gene expression level to a baseline gene expression level. The change in expression may be displayed as a graph on the display screen.

A further understanding of the nature and advantages of the inventions herein may be realized by reference to the remaining portions of the specification and the attached drawings.

## DESCRIPTION OF PREFERRED EMBODIMENTS

General

The present invention provides innovative methods of identifying nucleotides (*i.e.*, base calling) in sample nucleic acid sequences and monitoring gene expression. In the description that follows, the invention will be described in reference to preferred embodiments. However, the description is provided for purposes of illustration and not for limiting the spirit and scope of the invention.

Fig. 1 illustrates an example of a computer system that may be used to execute software embodiments of the present invention. Fig. 1 shows a computer system 1 which includes a monitor 3, screen 5, cabinet 7, keyboard 9, and mouse 11.

Mouse 11 may have one or more buttons such as mouse buttons 13. Cabinet 7 houses a CD-ROM drive 15 and a hard drive (not shown) that may be utilized to store and retrieve software programs including computer code incorporating the present invention. Although a CD-ROM 17 is shown as the computer readable medium, other computer readable media including floppy disks, DRAM, hard drives, flash memory, tape, and the like may be utilized. Cabinet 7 also houses familiar computer components (not shown) such as a processor, memory, and the like.

Fig. 2 shows a system block diagram of computer system 1 used to execute software embodiments of the present invention. As in Fig. 1, computer system 1 includes monitor 3 and keyboard 9. Computer system 1 further includes subsystems such as a central processor 50, system memory 52, I/O controller 54, display adapter 56, removable disk 58, fixed disk 60, network interface 62, and speaker 64. Removable disk 58 is representative of removable computer readable media like floppies, tape, CD-ROM, removable hard drive, flash memory, and the like. Fixed disk 60 is representative of an internal hard drive or the like. Other computer systems suitable for use with the present invention may include additional or fewer subsystems. For example, another computer system could include more than one processor 50 (*i.e.*, a multi-processor system) or memory cache.

Arrows such as 66 represent the system bus architecture of computer

整理番号 = P A 5 2 C 2 9 5

ページ (11)

system 1. However, these arrows are illustrative of any interconnection scheme serving to link the subsystems. For example, display adapter 56 may be connected to central processor 50 through a local bus or the system may include a memory cache. Computer system 1 shown in Fig. 2 is but an example of a computer system suitable for use with the present invention. Other configurations of subsystems suitable for use with the present invention will be readily apparent to one of ordinary skill in the art. In one embodiment, the computer system is a workstation from Sun Microsystems.

The VLSIPS™ technology provides methods of making very large arrays of oligonucleotide probes on very small chips. See U.S. Patent No. 5,143,854 and PCT patent publication Nos. WO 90/15070 and 92/10092, each of which is hereby incorporated by reference for all purposes. The oligonucleotide probes on the chip are used to detect complementary nucleic acid sequences in a sample nucleic acid of interest (the "target" nucleic acid).

The present invention provides methods of analyzing hybridization intensity files for a chip containing hybridized nucleic acid probes. In a representative embodiment, the files represent fluorescence data from a biological array, but the files may also represent other data such as radioactive intensity data. Therefore, the present invention is not limited to analyzing fluorescent measurements of hybridizations but may be readily utilized to analyze other measurements of hybridization.

For purposes of illustration, the present invention is described as being part of a computer system that designs a chip mask, synthesizes the probes on the chip, labels the nucleic acids, and scans the hybridized nucleic acid probes. Such a system is fully described in U.S. Patent Application No. 08/249,188 which is hereby incorporated by reference for all purposes. However, the present invention may be used separately from the overall system for analyzing data generated by such systems, such as at remote locations.

Fig. 3 illustrates a computerized system for forming and analyzing arrays of biological materials such as RNA or DNA. A computer 100 is used to design arrays of biological polymers such as RNA or DNA. The computer 100 may be, for example, an appropriately programmed IBM personal computer compatible



整理番号=P A 5 2 C 2 9 5

ページ (12)

running Windows NT including appropriate memory and a CPU as shown in Figs. 1 and 2. The computer system 100 obtains inputs from a user regarding characteristics of a gene of interest, and other inputs regarding the desired features of the array. Optionally, the computer system may obtain information regarding a specific genetic sequence of interest from an external or internal database 102 such as GenBank. The output of the computer system 100 is a set of chip design computer files 104 in the form of, for example, a switch matrix, as described in PCT application WO 92/10092, and other associated computer files.

The chip design files are provided to a system 106 that designs the lithographic masks used in the fabrication of arrays of molecules such as DNA. The system or process 106 may include the hardware necessary to manufacture masks 110 and also the necessary computer hardware and software 108 necessary to lay the mask patterns out on the mask in an efficient manner. As with the other features in Fig. 3, such equipment may or may not be located at the same physical site, but is shown together for ease of illustration in Fig. 3. The system 106 generates masks 110 or other synthesis patterns such as chrome-on-glass masks for use in the fabrication of polymer arrays.

The masks 110, as well as selected information relating to the design of the chips from system 100, are used in a synthesis system 112. Synthesis system 112 includes the necessary hardware and software used to fabricate arrays of polymers on a substrate or chip 114. For example, synthesizer 112 includes a light source 116 and a chemical flow cell 118 on which the substrate or chip 114 is placed.

Mask 110 is placed between the light source and the substrate/chip, and the two are translated relative to each other at appropriate times for deprotection of selected regions of the chip. Selected chemical reagents are directed through flow cell 118 for coupling to deprotected regions, as well as for washing and other operations. All operations are preferably directed by an appropriately programmed computer 119, which may or may not be the same computer as the computer(s) used in mask design and mask making.

The substrates fabricated by synthesis system 112 are optionally diced into smaller chips and exposed to marked targets. The targets may or may not be complementary to one or more of the molecules on the substrate. The targets are

整理番号 = P A 5 2 0 2 9 5

ページ (13)

marked with a label such as a fluorescein label (indicated by an asterisk in Fig. 3) and placed in scanning system 120. Scanning system 120 again operates under the direction of an appropriately programmed digital computer 122, which also may or may not be the same computer as the computers used in synthesis, mask making, and mask design. The scanner 120 includes a detection device 124 such as a confocal microscope or CCD (charge-coupled device) that is used to detect the location where labeled target (\*) has bound to the substrate. The output of scanner 120 is an image file(s) 124 indicating, in the case of fluorescein labeled target, the fluorescence intensity (photon counts or other related measurements, such as voltage) as a function of position on the substrate. Since higher photon counts will be observed where the labeled target has bound more strongly to the array of polymers, and since the monomer sequence of the polymers on the substrate is known as a function of position, it becomes possible to determine the sequence(s) of polymer(s) on the substrate that are complementary to the target.

The image file 124 is provided as input to an analysis system 126 that incorporates the visualization and analysis methods of the present invention. Again, the analysis system may be any one of a wide variety of computer system(s).

The present invention provides various methods of analyzing the chip design files and the image files, providing appropriate output 128. The present invention may further be used to identify specific mutations in a target such as DNA or RNA.

Fig. 4 provides a simplified illustration of the overall software system used in the operation of one embodiment of the invention. As shown in Fig. 4, the system first identifies the genetic sequence(s) or targets that would be of interest in a particular analysis at step 202. The sequences of interest may, for example, be normal or mutant portions of a gene, genes that identify heredity, or provide forensic information. Sequence selection may be provided via manual input of text files or may be from external sources such as GenBank. At step 204 the system evaluates the gene to determine or assist the user in determining which probes would be desirable on the chip, and provides an appropriate "layout" on the chip for the probes.

The chip usually includes probes that are complementary to a reference nucleic acid sequence which has a known sequence. A wild-type probe

整理番号 = P A 5 2 C 2 9 5

ページ (14)

is a probe that will ideally hybridize with the reference sequence and thus a wild-type gene (also called the chip wild-type) would ideally hybridize with wild-type probes on the chip. The target sequence is substantially similar to the reference sequence except for the presence of mutations, insertions, deletions, and the like. The layout implements desired characteristics such as arrangement on the chip that permits "reading" of genetic sequence and/or minimization of edge effects, ease of synthesis, and the like.

Fig. 5 illustrates the global layout of a chip. Chip 114 is composed of multiple units where each unit may contain different tilings for the wild-type sequence or multiple wild-type sequences. Unit 1 is shown in greater detail and shows that each unit is composed of multiple cells which are areas on the chip that may contain probes. Conceptually, each unit includes multiple sets of related cells.

As used herein, the term cell refers to a region on a substrate that contains many copies of a molecule or molecules (e.g., nucleic acid probes).

Each unit is composed of multiple cells that may be placed in rows (or "lanes") and columns. In one embodiment, a set of five related cells includes the following: a wild-type cell 220, "mutation" cells 222, and a "blank" cell 224. Cell 220 contains a wild-type probe that is the complement of a portion of the wild-type sequence. Cells 222 contain "mutation" probes for the wild-type sequence. For example, if the wild-type probe is 3'-ACGT, the probes 3'-ACAT, 3'-ACCT, 3'-ACGT, and 3'-ACTT may be the "mutation" probes. Cell 224 is the "blank" cell because it contains no probes (also called the "blank" probe). As the blank cell contains no probes, labeled targets should not bind to the chip in this area. Thus, the blank cell provides an area that can be used to measure the background intensity.

Again referring to Fig. 4, at step 206 the masks for the synthesis are designed. At step 208 the software utilizes the mask design and layout information to make the DNA or other polymer chips. This software 208 will control, among other things, relative translation of a substrate and the mask, the flow of desired reagents through a flow cell, the synthesis temperature of the flow cell, and other parameters. At step 210, another piece of software is used in scanning a chip thus synthesized and exposed to a labeled target. The software controls the scanning of the chip, and stores the data thus obtained in a file that may

整理番号 = P A 5 2 C 2 9 5

ページ (15)

later be utilized to extract sequence information.

At step 212 a computer system utilizes the layout information and the fluorescence information to evaluate the hybridized nucleic acid probes on the chip.

Among the important pieces of information obtained from DNA chips are the identification of mutant targets and determination of genetic sequence of a particular target.

Fig. 6 illustrates the binding of a particular target DNA to an array of DNA probes 114. As shown in this simple example, the following probes are formed in the array (only one probe is shown for the wild-type probe):

```

3'-AGAACGT
   AGACCGT
   AGAGCGT
   AGATCGT

```

```

.
.
.

```

As shown, the set of probes differ by only one base, a single base mismatch at an interrogation position, so the probes are designed to determine the identity of the base at that location in the nucleic acid sequence. Accordingly, when used herein a unit will refer to multiple sets of related probes, where each set includes probes that differ by a single base mismatch at an interrogation position.

When a fluorescein-labeled (or other marked) target with the sequence 5'-TCTTGCA is exposed to the array, it is complementary only to the probe 3'-AGAACGT, and fluorescein will be primarily found on the surface of the chip where 3'-AGAACGT is located. Thus, for each set of probes that differ by only one base, the image file will contain four fluorescence intensities, one for each probe. Each fluorescence intensity can therefore be associated with the nucleotide or base of each probe that is different from the other probes. Additionally, the image file will contain a "blank" cell which can be used as the fluorescence intensity of the background. By analyzing the five fluorescence intensities associated with a

整理番号 = P A 5 2 C 2 9 5

ページ (16)

specific base location, it becomes possible to extract sequence information from such arrays using the methods of the invention disclosed herein.

Fig. 7 illustrates probes arranged in lanes on a chip. A reference sequence (or chip wild-type sequence) is shown with five interrogation positions marked with number subscripts. An interrogation position is oftentimes a base position in the reference sequence where the target sequence may contain a mutation or otherwise differ from the reference sequence. The chip may contain five probe cells that correspond to each interrogation position. Each probe cell contains a set of probes that have a common base at the interrogation position. For example, at the first interrogation position,  $I_1$ , the reference sequence has a base T. The wild-type probe for this interrogation position is 3'-TGAC where the base A in the probe is complementary to the base at the interrogation position in the reference sequence.

Similarly, there are four "mutant" probe cells for the first interrogation position,  $I_1$ . The four mutant probes are 3'-TGAC, 3'-TGCC, 3'-TGGC, and 3'-TGTC. Each of the four mutant probes vary by a single base at the interrogation position. As shown, the wild-type and mutant probes are arranged in lanes on the chip. One of the mutant probes (in this case 3'-TGAC) is identical to the wild-type probe and therefore does not evidence a mutation. However, the redundancy gives a visual indication of mutations as will be seen in Fig. 8.

Still referring to Fig. 7, the chip contains wild-type and mutant probes for each of the other interrogation positions  $I_2$ - $I_5$ . In each case, the wild-type probe is equivalent to one of the mutant probes.

Fig. 8 illustrates a hybridization pattern of a target on a chip with a reference sequence as in Fig. 7. The reference sequence is shown along the top of the chip for comparison. The chip includes a WT-lane (wild-type), an A-lane, a C-lane, a G-lane, and a T-lane (or U). Each lane is a row of cells containing probes. The cells in the WT-lane contain probes that are complementary to the reference sequence. The cells in the A-, C-, G-, and T-lanes contain probes that are complementary to the reference sequence except that the named base is at the interrogation position.

In one embodiment, the hybridization of probes in a cell is determined

整理番号=PA520295

ページ (17)

by the fluorescent intensity (e.g., photon counts) of the cell resulting from the binding of marked target sequences. The fluorescent intensity may vary greatly among cells. For simplicity, Fig. 8 shows a high degree of hybridization by a cell containing a darkened area. The WT-lane allows a simple visual indication that there is a mutation at interrogation position  $I_4$  because the wild-type cell is not dark at that position. The cell in the C-lane is darkened which indicates that the mutation is from T->G (mutant probe cells are complementary so the C-cell indicates a G mutation). In a preferred embodiment, the WT-Lane is not utilized so four cells (not including any "blank" cell) are utilized to call a base at an interrogation position.

In practice, the fluorescent intensities of cells near an interrogation position having a mutation are relatively dark creating "dark regions" around a mutation. The lower fluorescent intensities result because the cells at interrogation positions near a mutation do not contain probes that are perfectly complementary to the target sequence; thus, the hybridization of these probes with the target sequence is lower. For example, the relative intensity of the cells at interrogation positions  $I_3$  and  $I_5$  may be relatively low because none of the probes therein are complementary to the target sequence. Although the lower fluorescent intensities reduce the resolution of the data, the methods of the present invention provide highly accurate base calling within the dark regions around a mutation and are able to identify other mutations within these regions.

Fig. 9 illustrates standard and alternate tilings on a chip. As shown, the chip includes twelve units (units<sub>1,12</sub>). Units<sub>1,4</sub> are tiled (i.e., designed and synthesized on the chip) to include probes complementary to the same reference sequence. For identification purposes, this group of units will be called the standard group. In general, base calls for the target sequence will be performed utilizing the standard group unless the invention determines that another group or groups should be utilized.

Units<sub>5,8</sub> are tiled to include probes complementary to the same reference sequence, but a reference sequence that differs from the reference sequence for the standard group. This group of units will be called an alternate group. Units<sub>9,12</sub> comprises another alternate group that are based on a reference

整理番号=PA520295

ページ (18)

sequence that is different from the reference sequences of the standard and first alternate groups. Although the reference sequences are different, they are often quite similar. For example, the reference sequences may be slightly different mutations of HIV. Embodiments of the present invention evaluate and utilize information from tilings based on reference sequences that would typically not be used in base calling the target sequence.

The units within a group may include identical probes, probes of different structure, probes from the same or different chips, and the like. For example, one unit may include 5-mer probes with the interrogation position at the third position in probes. Another unit may include 10-mer probes with an interrogation position at the sixth position. Additionally, these units may have been tiled on the same or different chips.

The expanded section at the bottom left portion of Fig. 9 illustrates that each block of a unit typically includes four cells, denoted A, C, G, and T. The base designations specify which base is at the interrogation position of each probe within the cell. Typically, there are hundreds or thousands of identical nucleic probes within each cell.

Although in preferred embodiments the cells may be arranged adjacent to each other in sequential order along the reference sequence, there is no requirement that the cells be in any particular location as long as the location on the chip is determinable. Additionally, although it may be beneficial to synthesize the different groups on a single chip for consistency of experiments, the methods of the present invention may be advantageously utilized with data from different tilings on different chips.

#### Analyzing Target Sequences

Fig. 10 shows a screen display of hybridization intensities from a chip.

During analysis, the system receives an image file including the scanned image of the hybridized chip. In a preferred embodiment, the image file shows fluorescent intensities and locations that labeled target nucleic acid sequences or fragments bound to the chip.

A screen display 260 utilizes the common windowing graphical user

整理番号= P A 5 2 C 2 9 5

ページ (19)

interface. The user may select to display the image file for inspection. After the user selects the image file to be displayed, a window 262 is displayed that includes the image file. The image file shown includes multiple rows of A-, C-, G-, and T-lanes.

As the user moves the cursor over the displayed image file, a status bar 264 indicates the X and Y position of the cursor and the fluorescent intensity at that position. Additionally, the user is able to utilize the pointing device to select a rectangular area of the image file in order to manipulate the sub-image. For example, the user may magnify the subimage so that the individual cells may be seen more clearly. Additionally, the user may adjust the contrast of the intensities to bring to light some differences in hybridization intensity that is not apparent at the current contrast setting.

Fig. 11 is a flowchart of a process of computing a base call from hybridization intensities of related probes. When used herein, "related probes" are probes that differ by a nucleotide base at an interrogation position. Although typically the probes are identical except at the interrogation position, the probes may differ at other base positions as well. Accordingly, the related probes differ by at least one base.

At step 302 the hybridization intensities of the four related probes are adjusted by subtracting the background or "blank" cell intensity. Preferably, if a hybridization intensity is then less than or equal to zero, the hybridization intensity is set equal to a small positive number to prevent division by zero or negative numbers in future calculations.

At step 304, the hybridization intensities are sorted by intensity. The highest intensity is then compared to a predetermined background difference cutoff at step 306. The background difference cutoff is a number that specifies the hybridization intensity the highest intensity probe must be over the background intensity in order to correctly call the unknown base. Thus, the background adjusted base intensity must be greater than the background difference cutoff or the unknown base is deemed to be not accurately callable.

If the highest hybridization intensity of the related probes is not greater than the background difference cutoff, the unknown base is assigned the



整理番号=P A 5 2 C 2 9 5

ページ (20)

code 'N' (insufficient intensity) as shown at step 308. Otherwise, the ratio of the highest hybridization intensity and second highest hybridization intensity is calculated as shown at step 310.

At step 312, the ratio calculated at step 310 is compared to a predetermined ratio cutoff. The ratio cutoff is a number that specifies the ratio required to identify the unknown base. In preferred embodiments, the ratio cutoff is 1.2. If the ratio is greater than the ratio cutoff, the unknown base is called according to the probe with the highest hybridization intensity. Typically, the base is called as the complement of the base at the interrogation position in the highest intensity probe as shown at step 314. Otherwise, the ratio of the second highest hybridization intensity and third highest hybridization intensity is calculated as shown at step 316.

At step 318, the ratio calculated at step 316 is compared to the ratio cutoff. If the ratio is greater than the ratio cutoff, the unknown base is called as being an ambiguity code specifying the complements of interrogation position bases of the highest hybridization intensity probe and the second highest hybridization probe as shown at step 320. Otherwise, the ratio of the third highest hybridization intensity and fourth highest hybridization intensity is calculated as shown at step 322.

At step 324, the ratio calculated at step 322 is compared to the ratio cutoff. If the ratio is greater than the ratio cutoff, the unknown base is called as being an ambiguity code specifying the complements of interrogation position bases of the highest, second highest and third highest hybridization intensity probes as shown at step 326. Otherwise, the unknown base is assigned the code 'X' (insufficient discrimination) as shown at step 328.

Fig. 12 is a flowchart of another process of computing a base call from hybridization intensities of related probes. The flowchart shown operates on hybridization intensities demonstrated by related probes; thus, a base call is made for the base in the target corresponding to the interrogation position in probes that differ by a single base mismatch at the interrogation position. At step 402, the system determines if there is one probe with the highest hybridization to the target sequence. If there is not, the base is called as an 'N' meaning ambiguous. For

整理番号 = P A 5 2 C 2 9 5

ページ (21)

example, if two probes have the same highest intensity (i.e., there is a tie), the base would be called as 'N'.

If there is a single probe that has the highest hybridization to the target, the base is called according to that probe at step 406. Since the probes are complementary to the target sequence, the base may be called as the complementary base (C/G, A/T) to the base at the interrogation position of the probe.

At step 408, the system determines if the base call is a mutant, meaning it is different than the base in the reference sequence. If the base call is not a mutant base call, the base call has been made. Otherwise, the system determines checks to make sure certain "mutant" conditions are met at step 410 or the base is called as 'N' at step 412.

Before describing the mutant conditions for one embodiment, it may be beneficial to give labels to the hybridization intensities of the related probes. For illustration purposes "HighInt" will refer to the highest hybridization intensity, "SecondInt" will refer to the second highest hybridization intensity, "ThirdInt" will refer to the third highest hybridization intensity, and "LowInt" will refer to the lowest highest hybridization intensity.

In one embodiment, the mutant conditions include three tests that must all be met to call the base a mutant. A first test is whether the difference between HighInt and SecondInt is greater than a difference cutoff. Thus, the system determines if HighInt - SecondInt is greater than a predefined value. This value should be chosen to allow mutant base calls only when the highest hybridization intensity is greater than the next highest hybridization intensity by a desired amount.

A second test is whether a first ratio is less than a first ratio cutoff. The first ratio is the following:

$$\frac{\text{SecondInt} - \sqrt{\text{ThirdInt} * \text{LowInt}}}{\text{HighInt} - \sqrt{\text{ThirdInt} * \text{LowInt}}}$$

$$\text{HighInt} - \sqrt{\text{ThirdInt} * \text{LowInt}}$$

The system determines if this first ratio is less than a predefined value. This value should be chosen to allow mutant base calls only when the highest hybridization intensity is a desired ratio greater than the next highest hybridization intensity even after the lowest two hybridization intensities are subtracted out.

A third test is whether a neighbor ratio is greater than a neighbor ratio cutoff. The neighbor ratio is the following:

$$\frac{\text{HighInt}_n}{\text{HighInt}_n - \sqrt{(\text{HighInt}_{n+1} * \text{HighInt}_{n-1})}}$$

where the subscript n designates values for the base position that is being called and n+1 and n-1 represent values for adjacent base positions. Thus, the system determines if the neighbor ratio is greater than a predefined value. This value should be chosen to allow mutant base calls only when the highest hybridization intensity is a desired ratio greater than the highest hybridization intensity with the adjacent highest hybridization intensities subtracted out.

Accordingly, in a preferred embodiment, only if all of the mutant conditions are met will the base be called a mutant base. This embodiment recognizes that mutations are fairly rare so a mutant base should only be called when there is a high likelihood that there has been a mutation. If the mutant conditions are not met, the base may be called as ambiguous or as the same as the reference sequence (which statistically may be the correct base call).

Although a preferred embodiment utilizes three mutant conditions, other embodiments may use a single mutant condition (e.g., one of the conditions described above). Other embodiments may utilize other base calling methods including the ones described in the U.S. Patent Applications previously incorporated by reference.

Fig. 13 is a flowchart of a process of calling bases in a group of units. As indicated earlier, a unit includes multiple sets of related cells, where the related cells include probes that differ by a single base at an interrogation position. In a typical embodiment, the system initially receives input on the hybridization intensities (e.g., from the image data file produced by a scanner that scans the hybridized chip) and the structure of the probes that correspond to the hybridization intensities. In preferred embodiments, the background intensity (e.g., intensity measured from "blank" cells or other areas of the chip without probes) are subtracted from the measured hybridization intensities. The background subtracted hybridization intensities may also be limited to have a minimum hybridization intensity of 1 (e.g., one photon count).

整理番号=PA52C295

ページ (23)

The hybridization intensity describes the extent of hybridization that was measured between a probe (or multiple copies of a probe) and the target sequence. As an example, the hybridization intensity may refer to the mean of the photon counts recorded from a cell, the photon counts resulting from fluorescein labeled target sequences that bound to probes in the cell.

At step 452, the system gets a base position in the target sequence to be called. The system then computes a base call for each unit of the group at step 454.

Therefore, the hybridization intensities for the related cells of each unit at the base position are analyzed. With this analysis (embodiments of which were described in more detail in reference to Figs. 11 and 12), the system computes a base call for each unit. Thus, if there are five units in the group, five base calls may be produced.

The system analyzes the base calls of the units of the group at step 456 in order to compute a base call for the group. In one embodiment, the system calls the base according to the base which is called most often by the units. For example, if there are five units and the following base calls were made for each unit:

'T' - three units

'G' - one unit

'N' - one unit

The base will be called a T since three out of five units agree. Ties may be broken by analyzing other factors like the highest average hybridization intensity of the unit or units that call each base in the tie. In a preferred embodiment, the invention utilizes the process described in Fig. 15.

At step 458, it is determined whether there is next base position to analyze. The present invention may be utilized to call all the bases of a target nucleic acid sequence so the process may, in effect, "walk" through the base positions. Additionally, the invention may be utilized to call only certain base positions (e.g., mutation positions) so the process may skip certain base positions altogether.

Fig. 14 is a flowchart of a process of calling bases for multiple groups of units. As shown in Fig. 9, there may be multiple groups on one or more chips that are to be analyzed. The multiple groups may be tiled according to different

整理番号=P A 5 2 C 2 9 5

ページ (24)

reference sequences; however, this does not mean that all of their hybridization information may not be utilized. Typically, it is assumed that the reference sequence for the standard group is expected to be the most identical to the target sequence. However, if one of the alternate groups is determined to be more identical (i.e., better for making a base call), then that group will be used to make the base call.

At step 502, the system computes base calls in the units of the standard and alternate groups. The base calling may be done as was described in reference to Fig. 13.

The system then computes a base call for each group of units at step 504. This may be accomplished by determining the base that is called most often by the units. Alternatively, the base call for the group may be determined utilizing the process which will be described in more detail in reference to Fig. 15.

After the system has determined a base call for each group of units (both the standard and alternate tilings), the system identifies a base position at step 506. The system then determines the best group of units for this base position to be utilized to make the base call. In general, selecting the best group may involve determining which reference sequence of the groups has the fewest mismatches with the target sequence near or in a window around the interrogation position. The group of units that has the fewest mismatches near the interrogation position may have the highest likelihood of producing the most accurate base call. An embodiment of selecting the best group will be described in more detail in reference to Fig. 16.

At step 510, the system calls the base at the identified base position according to the best group of units (i.e., utilizing the base call for the group that was computed at step 504). Once the base call has been made, the system determines if there is a next base position to perform a base call. If there is another base position to be called, the system proceeds to call that base position at step 506.

Fig. 15 is a flowchart of a process of calling a base for a group of units.

At step 602, the system determines if a majority of units call the same base at the specified base position. The majority is determined upon reference to only those units that call a base (e.g., do not call as ambiguous or 'N'). For example, assume

整理番号 = P A 5 2 0 2 9 5

ページ (25)

that there are seven units and the following base calls have been made for the units:

'G' - three units

'T' - one unit

'N' - four units

Since three out of four of the nonambiguous base calls are 'G', the system will initially call the base as a 'G' for the group of units. The base will be called as the majority base unless an exception rule applies at step 604.

The exception rules specify conditions which dictate what base call should be made for the group of units. These rules may include conditions that change a majority base call and may include conditions to deal with situations when there is not a base call that a majority of units call. In a preferred embodiment, the exception rules include tie breaking rules which analyze the hybridization intensity of neighboring probes (e.g., one unit calls one base and another unit calls a different base). Additionally, the exception rules specify that if three units call different bases with one of the calls being for the reference base, the system should call the base as the reference for the group of units. Other exception rules are described in the Appendix.

At step 606, the system determines if an exception rule applies. If an exception rule does apply, the rule is applied at step 608.

Fig. 16 is a flowchart of a process of selecting a best group of units for performing a base call. Selecting the best group involves determining which reference sequence of the groups has the fewest mismatches with the target sequence near the interrogation position. The group of units that has the fewest mismatches near the interrogation position may have the highest likelihood of producing the most accurate base call. The window around the interrogation position which is analyzed may be a set value or set according to the probe structure. For example, if the maximum distance that the probes for all the groups extend from the interrogation position is eight base positions to one side of the interrogation position and ten base positions to the other side of the interrogation position, the window may be set as including this range of base positions.

At step 702, the system calculates mismatch scores for the standard and alternate groups of units. The mismatch score is an indication of how many

整理番号=PA52C295

ページ (28)

mismatches a reference sequence appears to have with the target sequence. In order to determine a mismatch score, the system may only analyze base positions where at least two of the reference sequences differ. Thus, if all the reference sequences are identical at a base position, this base position may be skipped.

At each base position where at least two reference sequences differ, the system determines if the base call for a group (the base call indicating the likely base in the target sequence) at each of these positions differs from the corresponding base of the reference sequence. If the base call and the base for the reference sequence differ, the mismatch score is incremented by one. Initially, the mismatch scores for each group is set to zero.

Conceptually, it should be understood that the mismatch score is an indication of the number of base positions in a portion of the reference sequence that differ from the target sequence (optionally excluding those positions where all the reference sequences are the same). To better illustrate this concept, the following simple example is presented. Assume there is a standard group and two alternate groups as follows:

	<u>Standard Group</u>	<u>Mismatch Score</u>
reference	ACGGATGAGATACGA	1
base calls	ACTGATGAGATACGA	
	<u>Alternate Group 1</u>	<u>Mismatch Score</u>
reference	ACTGATGAGATACGA	0
base calls	ACTGATGAGATACGA	
	<u>Alternate Group 2</u>	<u>Mismatch Score</u>
reference	ACGGATGAGATACGT	2
base calls	ACTGATGAGATACGA	

The underlined bases correspond to the base position which is being analyzed. The bolded base positions indicate base positions where at least two of the reference sequences differ. At these bolded base positions, the standard group has one base position where the reference sequence differs from the target sequence (as indicated by the base calls) so the mismatch score is 1. Similarly, the first alternate group has

整理番号=P A 5 2 C 2 9 5

ページ (27)

a mismatch score of 0 and the second alternate group has a mismatch score of 2.

As alternate group 1 has the lowest mismatch score, that group would be utilized to call the base at the base position being analyzed. In this simple example, the base call is not different for any of the groups as this example is intended to illustrate how the best group may be selected. However, what is important is that the invention recognizes that the more mismatches that occur near a base position, the less accurate the base call will become. This result is brought upon by the fact that a mismatch between the reference sequence and the target sequence creates any area where the probes interrogating neighboring base positions include a single base mismatch. Single base mismatches lower the hybridization intensity and may produce inaccurate results.

At step 704, the system determines if a mismatch score of the standard groups is less than or equal to the mismatch scores of alternate groups. If the standard group has the lowest mismatch score (or ties), then the base call performed according to the standard group.

The system determines if a single alternate group has the lowest mismatch score at step 708. If so, that alternate group is utilized to make the base call at step 710. Otherwise, there are more than one alternate groups that have the same mismatch scores. If this is the case, the alternate group may be chosen which includes units that most consistently called the base at step 712. For example, if two alternate groups have the same lowest mismatch score but one group's units all called the same base and the other group's units were split, the alternate group that called the same base would be utilized. Other methods of determining the best group in the event of a mismatch score tie may also be utilized.

Fig. 17A shows a screen displays allowing analysis of nucleotides from experiments from one or more chips. A screen display 802 includes multiple screen areas that display different information to the user. A screen area 804 includes the name of a reference sequence which in this example is PRT 440A which are antisense regions (Protease Reverse Transcriptase) of the HIV virus. The reference sequence is typically used as a baseline with which to compare sample sequences. Although the reference sequence on the screen may be the chip wild-type sequence for which the chips were tiled, there is no requirement that this is the



整理番号= P A 5 2 C 2 9 5

ページ (28)

case.

A screen area 806 includes the nucleotide sequence for the reference sequence for the probe array. The base position of each nucleotide is shown above screen area 806. Screen area 806 also shows the reference sequence for each unit if "expanded" in the user interface.

A screen area 808 shows the user the chip and composite files that are currently being analyzed. A chip file (e.g., ends in ".CHP") includes data obtained from a single chip. A composite file (e.g., ends in ".CMP") includes data obtained from more than one chip. When a user opens a chip or composite file for analysis, the pathname of the file is displayed in screen area 808.

Information from the chip and composite files may be displayed in screen areas 810 and 812. Screen area 810 includes the names of sample sequences currently being analyzed from the chip or composite files. The name of the sample sequence is typically chosen to enable the user to readily determine the what the sample sequence represents. Screen area 812 includes the nucleotide sequence for the sample sequences. The base position of each nucleotide in screen area is the same as indicated above screen area 806. Accordingly, the system automatically aligns the reference and sample sequences for easier analysis.

Fig. 17A has been described in order to familiarize the reader with the layout of the screen display. However, as illustrated by Fig. 17B, the invention allows the user to hide (not display) and summarize information from chip and composite files. For example, if a user "clicks on" or activates the screen icon plus sign in front of the composite filename in screen area 808, the system displays more information about the composite file. As shown, the method that was utilized to combine the information from the chip files may be shown along with the individual chip files.

Additionally, if a user activates the screen icon plus sign in front of the chip filename in screen area 808, the system displays more information about the chip file including the process or procedure that was utilized to call bases. In Fig. 17B, the base calling procedure was the "Ratio Base Algorithm" which was described in reference to Fig. 10. Additionally, the user is able to modify parameters for the base calling procedure which will be immediately reflected in the

整理番号 = P A 5 2 C 2 9 5

ページ (29)

base calls shown on the display screen. For example, the ratio cutoff ("Ratio") is displayed as 1.2. If a user increases the ratio cutoff to 1.4, the system would then recalculate the base calls for the chip and the new base calls would be reflected in screen area 812. The parameters may be any values the are input into the base calling procedure including constants, thresholds, ranges, and the like.

Fig. 17B also illustrates that the system is able to combine data from multiple experiments (including various tilings) for easier reading of the user. The sample sequence 440-2A was shown in Fig. 17A and has been expanded in Fig. 17B to show that the base calls are derived from multiple experiments, where the data from multiple experiments may be from a single chip or multiple chips. In other words, the nucleotide sequence shown for sample sequence 440-2A in Figs. 17A and 17B do not represent a single experiment but actually a combination or consensus from multiple experiments. The user is able to review the data from each of the experiments as shown in Fig. 17B which includes displaying the hybridization intensities for each related base. The system allows the user to highlight a base position for analysis as shown for base position 100.

Referring again to Fig. 17A, a screen icon plus sign is displayed in front of the name of the sample sequence "440-2A." By activating the screen icon, the system displays each of the individual calls that make up the composite base call.

As shown in Fig. 17B, the composite base call is derived from multiple base calls. The multiple base calls are aligned with the composite base call according to base position. The invention provides great flexibility to the user for displaying, hiding, and summarizing data for analyzing sequences.

#### Monitoring Gene Expression

Fig. 18 shows a high level flowchart of a process of monitoring the expression of a gene by comparing hybridization intensities of pairs of perfect match and mismatch probes. The term "perfect match probe" refers to a probe that has a sequence that is perfectly complementary to a particular target sequence. The test probe is typically perfectly complementary to a portion (subsequence) of the target sequence. The term "mismatch control" or "mismatch probe" refer to probes whose sequence is deliberately selected not to be perfectly complementary to

整理番号 = P A 5 2 C 2 9 5

ページ (30)

a particular target sequence. For each mismatch (MM) control in a high-density array there typically exists a corresponding perfect match (PM) probe that is perfectly complementary to the same particular target sequence.

The process compares hybridization intensities of pairs of perfect match and mismatch probes that are preferably covalently attached to the surface of a substrate or chip. Most preferably, the nucleic acid probes have a density greater than about 60 different nucleic acid probes per  $1 \text{ cm}^2$  of the substrate. Although the flowcharts show a sequence of steps for clarity, this is not an indication that the steps must be performed in this specific order. One of ordinary skill in the art would readily recognize that many of the steps may be reordered, combined, and deleted without departing from the invention.

Initially, nucleic acid probes are selected that are complementary to the target sequence (or gene). These probes are the perfect match probes. Another set of probes is specified that are intended to be not perfectly complementary to the target sequence. These probes are the mismatch probes and each mismatch probe includes at least one nucleotide mismatch from a perfect match probe. Accordingly, a mismatch probe and the perfect match probe from which it was derived make up a pair of probes. As mentioned earlier, the nucleotide mismatch is preferably near the center of the mismatch probe.

The probe lengths of the perfect match probes are typically chosen to exhibit high hybridization affinity with the target sequence. For example, the nucleic acid probes may be all 20-mers. However, probes of varying lengths may also be synthesized on the substrate for any number of reasons including resolving ambiguities.

The target sequence is typically fragmented, labeled and exposed to a substrate including the nucleic acid probes as described earlier. The hybridization intensities of the nucleic acid probes is then measured and input into a computer system. The computer system may be the same system that directs the substrate hybridization or it may be a different system altogether. Of course, any computer system for use with the invention should have available other details of the experiment including possibly the gene name, gene sequence, probe sequences, probe locations on the substrate, and the like.

整理番号=PA520295

ページ (31)

Referring to Fig. 18, after hybridization, the computer system receives input of hybridization intensities of the multiple pairs of perfect match and mismatch probes at step 902. The hybridization intensities indicate hybridization affinity between the nucleic acid probes and the target nucleic acid (which corresponds to a gene). Each pair includes a perfect match probe that is perfectly complementary to a portion of the target nucleic acid and a mismatch probe that differs from the perfect match probe by at least one nucleotide.

At step 904, the computer system compares the hybridization intensities of the perfect match and mismatch probes of each pair. If the gene is expressed, the hybridization intensity (or affinity) of a perfect match probe of a pair should be recognizably higher than the corresponding mismatch probe. Generally, if the hybridizations intensities of a pair of probes are substantially the same, it may indicate the gene is not expressed. However, the determination is not based on a single pair of probes, the determination of whether a gene is expressed is based on an analysis of many pairs of probes. An exemplary process of comparing the hybridization intensities of the pairs of probes will be described in more detail in reference to Fig. 19.

After the system compares the hybridization intensity of the perfect match and mismatch probes, the system indicates expression of the gene at step 906.

As an example, the system may indicate an expression call to a user that the gene is either present (expressed), marginal or absent (unexpressed).

Fig. 19 shows a flowchart of a process of determining if a gene is expressed utilizing a decision matrix. At step 952, the computer system receives raw scan data of N pairs of perfect match and mismatch probes. In a preferred embodiment, the hybridization intensities are photon counts from a fluorescein labeled target that has hybridized to the probes on the substrate. For simplicity, the hybridization intensity of a perfect match probe will be designed " $I_{pm}$ " and the hybridization intensity of a mismatch probe will be designed " $I_{mis}$ ".

Hybridization intensities for a pair of probes is retrieved at step 954. The background signal intensity is subtracted from each of the hybridization intensities of the pair at step 956. Background subtraction may also be performed on all the raw scan data at the same time.

整理番号=P A 5 2 0 2 9 5

ページ (32)

At step 958, the hybridization intensities of the pair of probes are compared to a difference threshold (D) and a ratio threshold (R). It is determined if the difference between the hybridization intensities of the pair ( $I_{pm} - I_{nm}$ ) is greater than or equal to the difference threshold AND the quotient of the hybridization intensities of the pair ( $I_{pm} / I_{nm}$ ) is greater than or equal to the ratio threshold. The difference thresholds are typically user defined values that have been determined to produce accurate expression monitoring of a gene or genes. In one embodiment, the difference threshold is 20 and the ratio threshold is 1.2.

If  $I_{pm} - I_{nm} \geq D$  and  $I_{pm} / I_{nm} \geq R$ , the value NPOS is incremented at step 960. In general, NPOS is a value that indicates the number of pairs of probes which have hybridization intensities indicating that the gene is likely expressed. NPOS is utilized in a determination of the expression of the gene.

At step 962, it is determined if  $I_{nm} - I_{pm} \geq D$  and  $I_{nm} / I_{pm} \geq R$ . If this expression is true, the value NNEG is incremented at step 964. In general, NNEG is a value that indicates the number of pairs of probes which have hybridization intensities indicating that the gene is likely not expressed. NNEG, like NPOS, is utilized in a determination of the expression of the gene.

For each pair that exhibits hybridization intensities either indicating the gene is expressed or not expressed, a log ratio value (LR) and intensity difference value (IDIF) are calculated at step 966. LR is calculated by the log of the quotient of the hybridization intensities of the pair ( $I_{pm} / I_{nm}$ ). The IDIF is calculated by the difference between the hybridization intensities of the pair ( $I_{pm} - I_{nm}$ ). If there is a next pair of hybridization intensities at step 968, they are retrieved at step 954.

At step 972, a decision matrix is utilized to indicate if the gene is expressed. The decision matrix utilizes the values N, NPOS, NNEG, and LR (multiple LR's). The following four assignments are performed:

$$P1 = NPOS / NNEG$$

$$P2 = NPOS / N$$

$$P3 = (10 * \text{SUM}(\text{LR})) / (NPOS + NNEG)$$

These P values are then utilized to determine if the gene is expressed.

For purposes of illustration, the P values are broken down into ranges.

整理番号 = P A 5 2 C 2 9 5

ページ (33)

If P1 is greater than or equal to 2.1, then A is true. If P1 is less than 2.1 and greater than or equal to 1.8, then B is true. Otherwise, C is true. Thus, P1 is broken down into three ranges A, B and C. This is done to aid the readers understanding of the invention.

Thus, all of the P values are broken down into ranges according to the following:

$$A = (P1 \geq 2.1)$$

$$B = (2.1 > P1 \geq 1.8)$$

$$C = (P1 < 1.8)$$

$$X = (P2 \geq 0.35)$$

$$Y = (0.35 > P2 \geq 0.20)$$

$$Z = (P2 < 0.20)$$

$$Q = (P3 \geq 1.5)$$

$$R = (1.5 > P3 \geq 1.1)$$

$$S = (P3 < 1.1)$$

Once the P values are broken down into ranges according to the above boolean values, the gene expression is determined.

The gene expression is indicated as present (expressed), marginal or absent (not expressed). The gene is indicated as expressed if the following expression is true: A and (X or Y) and (Q or R). In other words, the gene is indicated as expressed if  $P1 \geq 2.1$ ,  $P2 \geq 0.20$  and  $P3 \geq 1.1$ . Additionally, the gene is indicated as expressed if the following expression is true: B and X and Q.

With the foregoing explanation, the following is a summary of the gene expression indications:

Present	A and (X or Y) and (Q or R)
	B and X and Q
Marginal	A and X and S
	B and X and R
	B and Y and (Q or R)

整理番号 = P A 5 2 C 2 9 5

ページ (34)

Absent

All others cases (e.g., any C combination)

In the output to the user, present may be indicated as "P," marginal as "M" and absent as "A" at step 974.

Once all the pairs of probes have been processed and the expression of the gene indicated, an average of ten times the LRs is computed at step 975. Additionally, an average of the IDIF values for the probes that incremented NPOS and NNEG is calculated, which may be utilized as an expression level. These values may be utilized for quantitative comparisons of this experiments with other experiments.

Quantitative measurements may be performed at step 976. For example, the current experiment may be compared to a previous experiment (e.g., utilizing values calculated at step 970). Additionally, the experiment may be compared to hybridization intensities of RNA (such as from bacteria) present in the biological sample in a known quantity. In this manner, one may verify the correctness of the gene expression indication or call, modify threshold values, or perform any number of modifications of the preceding.

For simplicity, Fig. 19 was described in reference to a single gene. However, the process may be utilized on multiple genes in a biological sample. Therefore, any discussion of the analysis of a single gene is not an indication that the process may not be extended to processing multiple genes.

Fig. 20 shows a screen display layout of gene expression monitoring software. A screen display 1000 is divided into two sections: a graphics display area 1002 and a data display area 1004. The graphics display area is for displaying graphs which will aid the user in interpreting the data. The data display area is for displaying the underlying data so the user may evaluate the underlying data for gene expression.

As will be shown in subsequent screen displays, the data display area is preferably organized in a table having rows and columns. Each column has a heading indicating the data that resides in the column. Each row represents data from a single experiment or combination of experiments for a gene. The term "experiment" is used herein to describe a process that created data. For example, a

整理番号=PA52C295

ページ (35)

single image file of a hybridized chip may produce many "experiments" for a number of genes. Additionally, experiments may refer to data obtained from different chips.

Fig. 21A shows a screen display illustrating the analysis of a selected gene. A screen display 1030 includes a graphics display area that illustrates with bar graphs the hybridization intensities of perfect match probes and mismatch probes at each base position in a selected gene. The gene selected is shown highlighted in a data display area 1034.

The data display area includes a number of column headings. The Experiment Name refers to a user-defined name for the experiment. The Gene Name is the name of the gene. The numbers Positive and Negative refer to the values NPOS and NNEG as described in reference to Fig. 19. The Pairs column indicates the number of perfect match and mismatch probe pairs that were utilized in the analysis of the gene. The Pos Fraction column indicates the fraction of probe pairs that were scored as positive (i.e., Positive/Pairs).

The Avg Ratio column indicates the average of  $I_{pm}/I_{mm}$  for all probes for a gene. The Log Avg column indicates the average of the  $\log(I_{pm}/I_{mm})$ . The PM Excess column indicates the number of perfect match probes that exhibit a hybridization intensity above a user defined threshold. The MM Excess indicates the number of mismatch probes that exhibit a hybridization intensity above a user defined threshold. Referring now to Fig. 21B, the Pos/Neg column indicates ratio of the Positive column to the Negative column ("Inf" is utilized if the Negative column includes a zero). The Avg Diff column indicates the average intensity difference for the gene. The average intensity difference was computed at step 975 of Fig. 19 (i.e.,  $\text{average}(\text{DIF})$ ).

The Abs Call column indicates the gene expression call for the experiment. The values in this column may be "P" for present, "M" for marginal and "A" for absent. The gene expression call for a preferred embodiment is described in more detail in reference to step 974 of Fig. 19.

As the user selects an experiment, the graphics display area displays graphs to aid the user in interpreting the data. A button bar 1034 enables the user to select which graph or graphs to display in the graphics display area.



整理番号 = P A 5 2 C 2 9 5

ページ (36)

Additionally, the user is able to sort the data in the display data are according to values in a selected column.

Fig. 22 shows another screen display illustrating the analysis of a selected gene. A screen display 1060 includes a graphics display area 1062 illustrating a graph of the ratio of the hybridization intensity of the perfect match probe to the mismatch probe at each base position. The x-axis is the base position and the y-axis is the ratio of hybridization intensities. The statistical ratio threshold is plotted on the graph, which in this example is 1.2. This graph may be utilized by the user to analyze how many probe pairs ( $I_{pm}/I_{mm}$ ) are above or below the threshold. The graph also includes the gene and experiment names.

Fig. 23 shows a screen display illustrating the comparison of experiments for selected genes. A screen display 1160 includes a graphics display area 1062 and a data display area 1164. The graphics display area includes a graph of the ratio of the hybridization intensity of the perfect match probe to the mismatch probe at each base position for each of the experiments/genes selected in the data display area. In a preferred embodiment, the experiment name, gene name, and data plot are a different color for each gene to allow the user to more easily see the differences between or among selected genes.

Fig. 24 shows another screen display illustrating the comparison of experiments for selected genes. A screen display 1200 includes a graphics display area 1202 illustrating the expression levels of genes selected in a data display area 1204. The graph of the expression levels of the selected genes is a bar graph. In a preferred embodiment, the expression level is defined as the average intensity difference (see average(IDIF) in Fig. 19). The graph also includes the gene and experiment names.

Fig. 25 shows another screen display illustrating the comparison of experiments for selected genes with multiple graphs in the graphics display area. A screen display 1230 includes a graphics display area 1232 depicting multiple graphs for analyzing the genes selected in a data display area 1234. An expression level graph 1236, an average intensity difference graph 1238 and a hybridization intensity graph 1240 are shown for the selected genes.

Figs. 26A and 26B show the flow of a process of determining the

整理番号 = P A 5 2 C 2 9 5

ページ (37)

expression of a gene by comparing baseline scan data and experimental scan data. For example, the baseline scan data may be from a biological sample where it is known the gene is expressed. Thus, this scan data may be compared to a different biological sample to determine if the gene is expressed. Additionally, it may be determined how the expression of a gene or genes changes over time in a biological organism. Accordingly, the term "baseline" means that it will be used as a point of reference.

At step 1302, the computer system receives raw scan data of N pairs of perfect match and mismatch probes from the baseline. The hybridization intensity of a perfect match probe from the baseline will be designated " $I_{pm}$ " and the hybridization intensity of a mismatch probe from the baseline will be designated " $I_{mm}$ ." The background signal intensity is subtracted from each of the hybridization intensities of the pairs of baseline scan data at step 1304.

At step 1306, the computer system receives raw scan data of N pairs of perfect match and mismatch probes from the experimental biological sample. The hybridization intensity of a perfect match probes from the experiment will be designated " $J_{pm}$ " and the hybridization intensity of a mismatch probe from the experiment will be designated " $J_{mm}$ ." The background signal intensity is subtracted from each of the hybridization intensities of the pairs of experimental scan data at step 1308.

The hybridization intensities of an I and J pair may be normalized at step 1310. For example, the hybridization intensities of the I and J pairs may be divided by the hybridization intensity of control probes.

At step 1312, the hybridization intensities of the I and J pair of probes are compared to a difference threshold (DDIF) and a ratio threshold (RDIF). It is determined if the difference between the hybridization intensities of the one pair ( $J_{pm} - J_{mm}$ ) and the other pair ( $I_{pm} - I_{mm}$ ) are greater than or equal to the difference threshold AND the quotient of the hybridization intensities of one pair ( $J_{pm} - J_{mm}$ ) and the other pair ( $I_{pm} - I_{mm}$ ) are greater than or equal to the ratio threshold. The difference thresholds are typically user defined values that have been determined to produce accurate expression monitoring of a gene or genes.

If  $(J_{pm} - J_{mm}) - (I_{pm} - I_{mm}) \geq DDIF$  and  $(J_{pm} - J_{mm}) / (I_{pm} - I_{mm}) \geq RDIF$ ,

整理番号 = P A 5 2 0 2 9 5

ページ (38)

the value NINC is incremented at step 1314. In general, NINC is a value that indicates the experimental pair of probes indicates that the gene expression is likely greater (or increased) than the baseline sample. NINC is utilized in a determination of whether the expression of the gene is greater (or increased), less (or decreased) or did not change in the experimental sample compared to the baseline sample.

At step 1316, it is determined if  $(J_{pm} - J_{bm}) - (I_{pb} - I_{bm}) \geq DDIF$  and  $(J_{pm} - J_{bm}) / (I_{pb} / I_{bm}) \geq RDIF$ . If this expression is true, NDEC is incremented. In general, NDEC is a value that indicates the experimental pair of probes indicates that the gene expression is likely less (or decreased) than the baseline sample. NDEC is utilized in a determination of whether the expression of the gene is greater (or increased), less (or decreased) or did not change in the experimental sample compared to the baseline sample.

For each of the pairs that exhibits hybridization intensities either indicating the gene is expressed more or less in the experimental sample, the values NPOS, NNEG and LR are calculated for each pair of probes. These values are calculated as discussed above in reference to Fig. 19. A suffix of either "B" or "E" has been added to each value in order to indicate if the value denotes the baseline sample or the experimental sample, respectively. If there are next pairs of hybridization intensities at step 1322, they are processed in a similar manner as shown.

Referring now to Fig. 26B, an absolute decision computation is performed for both the baseline and experimental samples at step 1324. The absolute decision computation is an indication of whether the gene is expressed, marginal or absent in each of the baseline and experimental samples. Accordingly, in a preferred embodiment, this step entails performing steps 972 and 974 from Fig. 19 for each of the samples. This being done, there is an indication of gene expression for each of the samples taken alone.

At step 1326, a decision matrix is utilized to determine the difference in gene expression between the two samples. This decision matrix utilizes the values, N, NPOSB, NPOSE, NNEGB, NNEGE, NINC, NDEC, LRB, and LRE as they were calculated above. The decision matrix performs different calculations

整理番号 = P A 5 2 C 2 9 5

ページ (39)

depending on whether NINC is greater than or equal to NDEC. The calculations are as follows.

If  $NINC \geq NDEC$ , the following four P values are determined:

$$P1 = NINC / NDEC$$

$$P2 = NINC / N$$

$$P3 = ((NPOSE - NPOSB) - (NNEGE - NNEGB)) / N$$

$$P4 = 10 * SUM(LRE - LRB) / N$$

These P values are then utilized to determine the difference in gene expression between the two samples.

For purposes of illustration, the P values are broken down into ranges as was done previously. Thus, all of the P values are broken down into ranges according to the following:

$$A = (P1 \geq 2.8)$$

$$B = (2.8 > P1 \geq 2.0)$$

$$C = (P1 < 2.0)$$

$$X = (P2 \geq 0.34)$$

$$Y = (0.34 > P2 \geq 0.24)$$

$$Z = (P2 < 0.24)$$

$$M = (P3 \geq 0.20)$$

$$N = (0.20 > P3 \geq 0.12)$$

$$O = (P3 < 0.12)$$

$$Q = (P4 \geq 0.9)$$

$$R = (0.9 > P4 \geq 0.5)$$

$$S = (P4 < 0.5)$$

Once the P values are broken down into ranges according to the above boolean values, the difference in gene expression between the two samples is determined.

In this case where  $NINC \geq NDEC$ , the gene expression change is indicated as increased, marginal increase or no change. The following is a

整理番号 = P A 5 2 C 2 9 5

ページ (40)

summary of the gene expression indications:

Increased	A and (X or Y) and (Q or R) and (M or N or O) A and (X or Y) and (Q or R or S) and (M or N) B and (X or Y) and (Q or R) and (M or N) A and X and (Q or R or S) and (M or N or O)
Marginal	A or Y or S or O
Increase	B and (X or Y) and (Q or R) and O B and (X or Y) and S and (M or N) C and (X or Y) and (Q or R) and (M or N)
No Change	All others cases (e.g., any Z combination)

In the output to the user, increased may be indicated as "I," marginal increase as "MI" and no change as "NC."

If  $NINC < NDEC$ , the following four P values are determined:

$$P1 = NDEC / NINC$$

$$P2 = NDEC / N$$

$$P3 = ((NNEGE - NNEGB) - (NPOSE - NPOSB)) / N$$

$$P4 = 10 * SUM(LRE - LRB) / N$$

These P values are then utilized to determine the difference in gene expression between the two samples.

The P values are broken down into the same ranges as for the other case where  $NINC \geq NDEC$ . Thus, P values in this case indicate the same ranges and will not be repeated for the sake of brevity. However, the ranges generally indicate different changes in the gene expression between the two samples as shown below.

In this case where  $NINC < NDEC$ , the gene expression change is indicated as decreased, marginal decrease or no change. The following is a summary of the gene expression indications:

整理番号 = P A 5 2 C 2 9 5

ページ (41)

Decreased	A and (X or Y) and (Q or R) and (M or N or O) A and (X or Y) and (Q or R or S) and (M or N) B and (X or Y) and (Q or R) and (M or N) A and X and (Q or R or S) and (M or N or O)
Marginal	A or Y or S or O
Decrease	B and (X or Y) and (Q or R) and O B and (X or Y) and S and (M or N) C and (X or Y) and (Q or R) and (M or N)
No Change	All others cases (e.g., any Z combination)

In the output to the user, decreased may be indicated as "D," marginal decrease as "MD" and no change as "NC."

The above has shown that the relative difference between the gene expression between a baseline sample and an experimental sample may be determined. An additional test may be performed that would change an I, MI, D, or MD (i.e., not NC) call to NC if the gene is indicated as expressed in both samples (e.g., from step 1324) and the following expressions are all true:

$$\text{Average(IDIFB)} \geq 200$$

$$\text{Average(IDIFE)} \geq 200$$

$$1.4 \geq \text{Average(IDIFE)} / \text{Average(IDIFB)} \geq 0.7$$

Thus, when a gene is expressed in both samples, a call of increased or decreased (whether marginal or not) will be changed to a no change call if the average intensity difference for each sample is relatively large or substantially the same for both samples. The IDIFB and IDIFE are calculated as the sum of all the IDIFs for each sample divided by N.

At step 1328, values for quantitative difference evaluation are calculated. An average of  $(I_{pm} - I_{um}) - (I_{pn} - I_{nn})$  for each of the pairs is calculated.

整理番号= P A 5 2 C 2 9 5

ページ (42)

Additionally, a quotient of the average of  $J_{pm} - J_{mm}$  and the average of  $I_{pm} - I_{mm}$  is calculated. These values may be utilized to compare the results with other experiments in step 1330.

Fig. 27A shows a screen display illustrating the monitoring of the change of gene expression between experiments. A screen display 1400 includes a graphics display area 1402 and a data display area 1404. A user begins the comparison of experiments for a gene by selecting two experiments for a gene. For simplicity, we will call one baseline data and the other experimental data, meaning it will be compared to the baseline. For example, a user may select two experiments for the gene with the name "g102506." A comparison of two experiments is an experiment itself so the user is able to enter an experiment name which was entered as "foo" in the data display area of Fig. 27A. Fig. 27B shows another screen display illustrating monitoring of the change of gene expression between experiments.

The system then determines the change in gene expression between the selected experiments according to the process described in Figs. 28A and 28B. The data display area includes columns denoting the data produced by this comparison. The Experiment Name refers to a user-defined name for the comparison experiment. The Gene Name is the name of the gene. The numbers Inc and Dec refer to the values NINC and NDEC as described in reference to Fig. 26A. More specifically, Inc refers to the number of base positions in the gene for which the difference and ratio of the perfect match and mismatch hybridization intensities are significantly greater in the experimental data.

The Inc Ratio column indicates the number of base positions where the hybridization intensity increased divided by the total number of base positions in the gene which are analyzed. The Dec Ratio column indicates the number of base positions where the hybridization intensity decreased divided by the total number of base positions in the gene which are analyzed. The Pos Change column indicates the difference in the number of positive scoring probe pairs in the experimental data versus the baseline data. The Neg Change column indicates the difference in the number of negative scoring probe pairs (perfect match and mismatch) in the experimental data versus the baseline data.

整理番号=PA520295

ページ (43)

The Inc/Dec column indicates the number probe pairs which had an increase in hybridization intensity in the experimental data versus the number of probe pairs which had a decrease in hybridization intensity in the experimental data.

The Avg Diff column indicates the average intensity difference in the experimental data.

The Diff Call column (not shown) indicates the change in expression level between the experiments for the gene. The column shows a "I" for increased gene expression, "MI" for marginal increased gene expression, "D" for decreased gene expression, "MD" for marginal decreased gene expression, "NC" for no change, and "?" for unknown. In a preferred embodiment, the change in expression level is calculated as described in reference to step 1326 of Fig. 26B.

In addition to calculating the change in gene expression, the user may also select graphs to analyze the data. Graphics display area 1402 shows three different graphs depicting the data from the baseline and experimental data.

Fig. 28 shows a screen display illustrating a three-dimensional bar graph which illustrates the change of gene expression between experiments. A screen display 1440 displays a graphical display area 1442 including a three-dimensional bar graph of the expression level of selected genes in a data display area 1444. The user selects one or more genes in the data display area and then instructs the system to generate a three-dimensional bar graph of the expression level of these genes, where the expression level in a preferred embodiment is the average intensity difference (i.e., average(IDIF)). The three-dimensional bar graph allows the user to easily view the expression level of multiple genes. Additionally, similar genes selected from multiple experiments may be shown simultaneously and rotated to display differences in expression levels.

### Conclusion

The above description is illustrative and not restrictive. Many variations of the invention will become apparent to those of skill in the art upon review of this disclosure. Merely by way of example, while the invention is illustrated with particular reference to the evaluation of DNA (natural or unnatural), the methods can be used in the analysis from chips with other materials synthesized



整理番号 = P A 5 2 C 2 9 5

ページ (44)

thereon, such as RNA. The scope of the invention should, therefore, be determined not with reference to the above description, but instead should be determined with reference to the appended claims along with their full scope of equivalents.

#### 4 Brief Description of Drawings

Fig. 1 illustrates an example of a computer system that may be used to execute software embodiments of the present invention;

Fig. 2 shows a system block diagram of a typical computer system;

Fig. 3 illustrates an overall system for forming and analyzing arrays of biological materials such as DNA or RNA;

Fig. 4 is an illustration of an embodiment of software for the overall system;

Fig. 5 illustrates the global layout of a chip formed in the overall system;

Fig. 6 illustrates conceptually the binding of nucleic acid probes on chips to a labeled target;

Fig. 7 illustrates nucleic acid probes arranged in lanes on a chip;

Fig. 8 illustrates a hybridization pattern of a target on a chip with a reference sequence as in Fig. 7;

Fig. 9 illustrates standard and alternate tilings;

Fig. 10 shows a screen display of hybridization intensities from a chip;

Fig. 11 is a flowchart of a process of computing a base call from hybridization intensities of related probes;

Fig. 12 is a flowchart of another process of computing a base call from hybridization intensities of related probes;

Fig. 13 is a flowchart of a process of calling bases in a group of units;

Fig. 14 is a flowchart of a process of calling bases for multiple groups of units;

Fig. 15 is a flowchart of a process of calling a base for a group of units;

Fig. 16 is a flowchart of a process of selecting a best group of units for performing a base call;

Figs. 17A and 17B show screen displays allowing analysis of

整理番号 = P A 5 2 C 2 9 5

ページ (45)

nucleotides from experiments from one or more chips;

Fig. 18 shows a high level flowchart of a process of monitoring the expression of a gene by comparing hybridization intensities of pairs of perfect match and mismatch probes;

Fig. 19 shows a flowchart of a process of determining if a gene is expressed utilizing a decision matrix;

Fig. 20 shows a screen display layout of gene expression monitoring software;

Figs. 21A and 21B show screen displays illustrating the analysis of a selected gene;

Fig. 22 shows another screen display illustrating the analysis of a selected gene;

Fig. 23 shows a screen display illustrating the comparison of experiments for selected genes;

Fig. 24 shows another screen display illustrating the comparison of experiments for selected genes;

Fig. 25 shows another screen display illustrating the comparison of experiments for selected genes with multiple graphs in the graphics display area;

Figs. 26A and 26B show a flowchart of a process of determining the expression of a gene by comparing baseline scan data and experimental scan data;

Figs. 27A and 27B show screen displays illustrating the monitoring of the change of gene expression between experiments; and

Fig. 28 shows a screen display illustrating a three-dimensional bar graph which illustrates the change of gene expression between experiments.